

Crowd4SDG

Citizen Science for the Sustainable Development Goals

Deliverable 5.1

Initial report on relevance and quality-related considerations of citizen-science generated data

Deliverable identifier: D5.1 Due date: 30/04/2021 Justification for delay: due to exceptional activity as agreed with the Project Officer Document release date: 24/06/2021 Nature: Report Dissemination Level: Public Work Package: 5 Lead Beneficiary: UNITAR Contributing Beneficiaries: UNIGE

Document status: Final

Abstract:

Deliverable 5.1 presents the results of the research study that was conducted under Crowd4SDG project to improve the understanding of the perspectives of National Statistical Offices (NSOs) and policymakers on the potential and limitations related to the use of citizen science data for monitoring and reporting on SDGs. It includes summary findings, a number of case studies from the countries and international organizations (IOs) that have been experimenting with citizen science or citizen generated data and a list of recommendations for NSOs and IOs on how to leverage this type of data better. It also includes a proposed criteria framework for citizen science data but also other non-official data sources drawing on some emerging examples from countries as well as UN Fundamental Principles of Official Statistics and UN Quality Assurance Framework for Official Statistics.

For more information on Crowd4SDG, please check: http://www.crowd4sdg.eu/



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No. 872944.



Document history

	Name	Partner	Date
Authored by	Elena Proden	UNITAR	09/04/2021
Edited by	Elena Proden, Madina Imaralieva	UNITAR	29/04/2021
Reviewed by	Barbara Pernici Jose Luis Fernandez Marquez WP5 Sub-Advisory Group members	POLIMI UNIGE	29/04/2021
Revised by	Elena Proden	UNITAR	14/06/2021
Approved by	François Grey, Jose Luis Fernandez Marquez	UNIGE	24/06/2021



Table of Contents

Document history	2
Project Partners	4
Crowd4SDG in Brief	5
Grant Agreement description of the deliverable	6
1. Summary	7
2. Purpose and scope of the deliverable	9
3. Methodology	13
4. Official Statistics: principles and quality assurance	15
5. Main survey findings	18
5.1 Track A. Official Statistics Community	18
5.2 Track B. Policy-makers	26
6. Examples of citizen science data for SDGs and lessons learnt	32
6.1 Key findings from the interviews	32
6.2 Citizen generated data to monitor selected SDGs	34
6.3 Citizen-generated data to provide key statistics relevant for the SDGs	37
6.4 Citizen-generated data and local level / social monitoring of the SDGs	38
6.5 Citizen-generated data beyond the SDGs	39
7. Guidelines for citizen data producers	43
8. Overall recommendations	48
8.1 Strengthening capacities of NSOs to leverage citizen science data	48
8.2 Possible criteria to be included in guidelines / protocols for CGD	49
9. Comparative analysis of data quality approaches between official statistics and academia/citizen science and recommendations	52
10. Conclusions and outlook	54
11. References	55
Annex: list of abbreviations	57



Project Partners

	Partner name	Acronym	Country
1 (COO)	Université de Genève	UNIGE	СН
2	European Organization for Nuclear Research	CERN	СН
3	Agencia Estatal Consejo Superior de Investigaciones Científicas	CSIC	ES
4	Politecnico di Milano	POLIMI	IT
5	United Nations Institute for Training and Research	UNITAR	СН
6	Université de Paris	UP	FR















Crowd4SDG in Brief

The 17 Sustainable Development Goals (SDGs), launched by the UN in 2015, are underpinned by over 160 concrete targets and over 230 measurable indicators. Some of these indicators initially had no established measurement methodology. For others, many countries do not have the data collection capacity. Measuring progress towards the SDGs is thus a challenge for most national statistical offices.

The goal of the Crowd4SDG project is to research the extent to which Citizen Science (CS) can provide an essential source of non-traditional data for tracking progress towards the SDGs, as well as the ability of CS to generate social innovations that enable such progress. Based on shared expertise in crowdsourcing for disaster response, the transdisciplinary Crowd4SDG consortium of six partners is focusing on SDG 13, Climate Action, to explore new ways of applying CS for monitoring the impacts of extreme climate events and strengthening the resilience of communities to climate related disasters.

To achieve this goal, Crowd4SDG is initiating research on the applications of artificial intelligence and machine learning to enhance CS and explore the use of social media and other non-traditional data sources for more effective monitoring of SDGs by citizens. Crowd4SDG is using direct channels through consortium partner UNITAR to provide National Statistical Offices (NSOs) with recommendations on best practices for generating and exploiting CS data for tracking the SDGs.

To this end, Crowd4SDG rigorously assesses the quality of the scientific knowledge and usefulness of practical innovations occurring when teams develop new CS projects focusing on climate action. This occurs through three annual challenge-based innovation events, involving online and in-person coaching. A wide range of stakeholders, from the UN, governments, the private sector, NGOs, academia, innovation incubators and maker spaces are involved in advising the project and exploiting the scientific knowledge and technical innovations that it generates.

Crowd4SDG has six work packages. Besides Project Management (UNIGE) and Dissemination & Outreach (CERN), the project features work packages on: Enhancing CS Tools (CSIC, POLIMI) with AI and social media analysis features, to improve data quality and deliberation processes in CS; New Metrics for CS (UP), to track and improve innovation in CS project coaching events; Impact Assessment of CS (UNITAR) with a focus on the requirements of NSOs as end-users of CS data for SDG monitoring. At the core of the project is Project Deployment (UNIGE) based on a novel innovation cycle called GEAR (Gather, Evaluate, Accelerate, Refine), which runs once a year.

The GEAR cycles involve online selection and coaching of citizen-generated ideas for climate action, using the UNIGE Open Seventeen Challenge (O17). The most promising projects are accelerated during a two-week in-person Challenge-Based Innovation (CBI) course. Top projects receive further support at annual SDG conferences hosted at partner sites. GEAR cycles focus on specific aspects of Climate Action connected with other SDGs like Gender Equality.



Grant Agreement description of the deliverable

The Deliverable 5.1 "Report on relevance and quality-related considerations of citizenscience generated data" is produced under Task 5.1: Analysis on the relevance and qualityrelated considerations of CS projects data for SDGs (UNITAR). This task is led by UNITAR with contributions from UNIGE.

T5.1: The focus of this task will be consultations and a study on the perceived relevance and quality considerations of CS projects data and non-traditional data sources, drawing on social media data among other, to support national governments in monitoring progress and implementing the Agenda 2030. The global SDG indicator framework currently composed of 232 single indicators was negotiated and being further refined and developed by the Inter-Agency Expert Group on SDG indicators composed of the representatives of national statistical offices. It will be subject to a more in-depth review by UN Statistical Commission in 2020 and 2025. At the national level, countries compile their national SDG indicators drawing on the global indicator framework while at the same time aligning national indicators with national priorities. It is expected that they will use to the extent possible international standards in this work. The first component of this task will involve consultations regarding the relevance aspects to be conducted using the methodologies, results, and findings from WP2 related to the application of citizen science methods for generating data on climate hazards and cities, climate resilience and gender, and climate adaptation and rights. A number of criteria will be explored to define the usefulness of the data degenerated by citizen science as a source of information to monitor performance on the SDGs and where relevant serve as secondary source in the process of production of official statistics. Among potential advantages that these data can offer are: - availability of timely information which may be a critical factor for certain decisions, incl. DRR context and other policy areas that may require rapid response; - provision of information for proxy indicators where the methodology does not exist (this is particularly relevant for Tier 3 global SDG indicators but can also be relevant in certain countries for Tier 2 indicators where National Statistical Systems may be struggling with producing data on indicators with established methodologies due to low capacity, high costs or other factors); - improved granularity which may offer insights into specific situations of various population groups, in particular most disadvantaged; - availability of gualitative information. The second component of this task involves a comparative analysis of principles guiding quality requirements applied by NSOs and the quality assessment frameworks developed by the scientific community. These will provide a basis for producing quality-related considerations on citizen-generated data and feedback into the work of Task 2.3. The resulting conclusions and recommendations will be disseminated in order to help inform any future standardsetting work that may be undertaken by the UN or by National Statistical Offices in this area. The results of both components will be used as reference for the enhancement of CS tools in WP2 and the coaching provided through the GEAR (Gather, Evaluate, Accelerate and Refine) methodology.



1. Summary

One of the main objectives of this report and the study that has informed it was to help bridging the citizen science community, on the one hand, and official statistics and decisionmaker communities, on the other, by examining the needs and the quality requirements for data generated by citizen science projects in the light of significant needs related to monitoring and reporting progress on SDGs.

There are multiple definitions of citizen science as noted in a recent study by Haklay, M. et al. (2020)¹. In this report, citizen science is studied more specifically from the perspective of its potential to contribute to the production of data relevant for the design and implementation of public policies at national and local levels. The working definition of citizen science data will cover therefore data produced with contributions of citizens who choose to voluntarily contribute their time, knowledge, skills and/or data to help produce needed evidence, strengthen accountability, or develop locally-rooted solutions. This definition will incorporate what is often known to National Statistical Offices (NSOs) as citizen-generated data (CGD).

Overall, citizen science data are perceived as having high potential for monitoring specifically environmental indicators in developed countries. In developing countries, these data sources are deemed as complementary data sources to help fill in important data gaps in the traditional areas of statistical activities where lack of capacities, resources or other limitations prevent National Statistical Systems (NSSs) from covering those through household surveys. The potential of citizen science data is also acknowledged on some sensitive indicators where inputs from civil society, human rights institutions and citizens can help complete the picture or else in humanitarian contexts.

The number of countries experimenting with citizen science data remains small. However, an interest in leveraging this data source grows and there is a number of very promising examples of specific projects/cases. In addition, a number of countries started working on protocols or quality assurance procedures for such data. The lack of know-how remains a major limitation that this report aims to address.

The study has elucidated different perceptions related to the use of citizen science and citizen generated data in support of the Sustainable Development Goals (SDGs) across different countries and National Statistical Systems (NSSs).

Broadly, one can distinguish 3-4 different set-ups in which the use of citizen science and citizen generated data are considered by National Statistical Systems. First, in some National Statistical Offices, Units in charge of SDG monitoring have the mandate to explore the use of non-traditional data sources to complement the reporting on the SDG indicators often on their open data platforms. Second, in other countries, other Units – sector statistics or partnerships – are involved in similar projects with leveraging citizen generated data in their area of SDG monitoring. Some countries have been using citizen generated data or citizen science data to fill in data gaps in important sectors they could not cover through household surveys and administrative data (such as informal economy).

Third, citizen generated data are often used for local level monitoring – sometimes in combination with specific accountability and feedback mechanisms linking back to the service provider.

¹ Haklay, Muki, et al. "Contours of Citizen Science: A Vignette Study", Royal Society Open Science, 2020.



Fourth, citizen generated and citizen science data are used to provide new insights – in combination with data science and machine learning – on our societies in close-to-real time mode. Those projects are often run by experimental statistics or data science departments and published not as official but as experimental statistics.

These arrangements do not seem in principle to be mutually exclusive but given the limited efforts so far in most cases it was just of them.

This report first presents the objectives of the study. It then introduces the methodology. This is followed by a discussion of the official statistics principles and quality assurance approach. A detailed analysis of the survey results for 3 tracks of respondents is presented in the subsequent section. This is followed by the presentation of the case studies and lessons learnt from those countries that have experimented with citizen science data. It then provides an overview of the guidelines developed by NSOs and international organizations for citizen science data. The final section is dedicated to the recommendations for NSOs and international organizations on how to leverage data more effectively using both "active" and "passive" approaches and includes a list of quality assurance criteria that could be considered by NSOs using "passive" approaches to leveraging citizen science data.



2. Purpose and scope of the deliverable

The report presents the results of a study undertaken by UNITAR under the Crowd4SDG project in 2020-2021 in consultation with the consortium of Crowd4SDG partners comprising University of Geneva, Polytechnic University of Milan, CSIC (Consejo Superior de Investigaciones Científicas), University of Paris, UNITAR and CERN and in collaboration with United Nations Statistics Division (UNSD), Global Partnership for Sustainable Development Data (GPSDD), International Institute of Applied System Analysis (IIASA), UN Women, United Nations Convention on Biological Diversity (UNCBD) and UN Conference on Trade and Development (UNCTAD) who all run various projects and activities related to leveraging citizen science data for monitoring progress on the SDGs and provided valuable advice in preparation of the study and during report finalization as a sub-advisory group. More specifically, the study has been designed to produce Deliverable 5.1 under the Working Package 5 "Impact Assessment" focusing on bridging the citizen science community, on the one hand, and official statistics and decision-maker communities, on the other, by examining the needs and the quality requirements for data generated by citizen science projects. This WP5 will also look, in the last seven months of the project, at the impact of the Crowd4SDG project and result in the preparation of a report with an analysis of the results achieved on the utilization of citizen science in support of monitoring the SDGs.

The report includes the analysis of main findings, good practices and potential areas where national data needs could be addressed with the support of citizen science data providing a series of success stories. It also proposes a set of recommendations on how to ensure that citizen science data could be leveraged at a larger scale for monitoring and supporting the implementation of the SDGs at global, national and local levels. A special emphasis is placed on recommendations covering quality related aspects.

As part of the study design, some important definitional issues needed to be addressed first because there is no uniform approach in the analysis of the use or potential use of citizen science data for monitoring SDGs.

There are multiple definitions of citizen science as noted in a recent study by Haklay, M., et al. (2020) that uses a participatory approach to examine the notion and "demonstrates the plurality of understanding of what citizen science is and calls for an open understanding of what activities are included in the field". It underscores the legacy of two traditions defining our today's perceptions: that of the meeting point between public participation and knowledge production (Irwin A., 2015) and that of citizens scientists being volunteers contributing to field observations (Bonney R., et al. 2009)².

In this report, citizen science is studied more specifically from the perspective of its potential to contribute to the production of data relevant for the design and implementation of public policies at national and local levels but not limited to the field observations or data collection.

Representatives of the official statistics community may be more familiar with the term "citizen generated data" often associated with community data and monitoring. The Global Partnership for Sustainable Development Data (GPSDD) guide³ defines the latter in a narrower sense closer to that of citizen science data. In other words, various forms of big data that may sometimes be considered citizen-generated data (digital trace) are not

² Ibid, pp. 1-2.

³ "Choosing and engaging with Citizen-Generated Data. Guide", GPSDD, Open Knowledge International and Public Data Lab, 2018.



considered as such neither in GPSDD's guide nor in this study. Those are referred to as social media, mobile phone or bank transactions data as distinct categories. Some of the key features of the citizen science data as used in this study are the <u>deliberate and voluntary</u> <u>nature</u> of the contributions of the community members to the process of data generation. Citizen science in general, understood more narrowly, also implies <u>scientific validation</u> of the knowledge generated by citizen scientists. Here, however, citizen science data will not include this condition as the scientific validation can in some cases be replaced by the <u>NSO</u> <u>validation</u> where no academic institution is part of the partnership.

In this report, the working definition of citizen science data will cover therefore the data produced with contributions of citizens who choose to voluntarily contribute their time, knowledge, skills and/or their data to help produce needed evidence, strengthen accountability, or develop locally-rooted solutions. This would <u>include initiatives in which</u> citizens may be helping collect or analyse data (including social media data), contribute their knowledge (including local and traditional knowledge) to data production, report voluntarily data on themselves, or else consent to share their data for a scientific experiment or data as a public good. This <u>would exclude</u> purely Big Data projects where citizens have no specific intention to contribute their data – such as from digital trace - to data production as public good.

Box 1. Report definition of citizen science data

The working definition of citizen science data will cover therefore the data produced with contributions of citizens who choose to voluntarily contribute their time, knowledge, skills and/or their data to help produce needed evidence, strengthen accountability, or develop locally-rooted solutions.

This definition will incorporate what is often known as citizen-generated data (CGD) but not the digital trace alone. The term "citizen generated data" will be used in some of the case studies where countries or international organizations used it to describe their activities. It will often match the definition of citizen science data used in this report. This approach also matches the approach proposed in the 2020 IIASA study (Fraisl, et al. (2020))⁴ that used the concept of "Public Participation in Scientific Research (PPSR)" from Shirk et al. (2012). This concept describes citizen science data as covering all of the below:

- Community-based monitoring (environmental monitoring and adaptive management activities by citizens in local communities;
- <u>Community-based participatory research</u> (knowledge co-creation with researchers, practitioners and community members working together to address issues relevant to communities, often those with a history of marginalization;
- Participatory action research (qualitative research methodology where communities are involved in all stages of defining the research process and with social change as main goal that played already an important role in the area of education;
- <u>Citizen-generated data sets</u> (data that people or their organisations produce to directly monitor, demand or drive change on issues that affect them, actively given by citizens, providing direct representations of their perspectives;
- <u>Crowdsourcing</u> (outsourcing the function to a large, undefined group of individuals through an open call;

⁴ Fraisl, D., et al. (2020). Mapping citizen science contributions to the UN sustainable development goals.



- <u>Volunteered geographic information</u> (digital geographical information that is generated and shared by individuals;
- <u>Participatory</u> sensing (an open practice of data capture, analysis and sharing through digital devices and platforms.

The main rationale behind this work and the report was a growing awareness among official statistics community that citizen science data offer a number of opportunities to help address the issues of <u>timeliness</u> and <u>granularity of data</u> for decision-making. It may often serve as a <u>useful complement</u> in some areas to official statistics, particularly, when it comes to global and national SDG indicators. It can also be <u>a standalone source in a number of areas</u>, often at a community level.

The presented recommendations are expected to address key data quality issues as well as other relevant aspects such as impediments and enabling factors on the side of the official statistics community, decision-makers and citizen science communities.

A number of organizations has already done some work to promote the use of citizen science data to support the monitoring of the SDGs. In particular, a paper was published in July 2020 by a group of scholars led by IIASA that has also included contributions from UNEP, CSCZ, and other experts. This paper has mapped all the global SDG indicators by examining the potential of using citizen science data for their production. The study has shown that 5 global indicators - mostly related to environment - are already produced with support from citizen science data, and 76 global indicators more could benefit from citizen science data (Fraisl, et al. 2020) ⁵. The awareness around citizen science data among National Statistical Offices, National Governments and International Organizations remains low, and it is not used to its full potential. IIASA has initiated a national pilot in Ghana supporting the Ghana Statistical Services (GSS) in addressing their data gaps and needs using citizen science data for several SDG indicators, including 14.1.1b. Floating Plastic Debris Density identified by the GSS based on the results of the systematic review undertaken by Fraisl, et al. (2020). In parallel, the UNDESA Statistics Division, in collaboration with GPSDD, the World Bank and the Sustainable Development Solutions Network (SDSN) have launched the Data4Now initiative with Colombia, Paraguay, Rwanda, Ghana, Senegal, Mongolia, Nepal, Bangladesh as pilot countries where citizen science data could be examined as one of the potential data sources. GPSDD have also developed guidelines on how to run CGD projects. While guidelines exist on how to run citizen science projects, the specificity of projects that could be run to provide data for NSOs is still to be further assessed and a set of quality and needs related recommendations would be useful.

Under CROWD4SDG, a two-pronged approach was proposed to promote the use of citizen science data for monitoring SDGs. First, the current study aims to assess the awareness and perceptions around citizen science data among official statistics communities and decision-makers, particularly, at the national level, as well as provide a stock-take of the lessons learnt from citizen science/ CGD projects where available or from the efforts to use other non-traditional data sources, with a special focus on quality requirements and needs. Second, depending on countries' priorities in the Data4Now initiative selecting this as a data source, we will jointly aim with UNSD to assess data coming from citizen science projects. In terms of topics, all SDGs were covered in the survey, but special attention was paid to survey results demonstrating an interest in topics relevant for the Crowd4SDG project such as climate change and resilience, sustainable human settlements, gender, water, Goal 16 or else biodiversity.

⁵ Ibid.



Based on this report, UNITAR will work with UNSD, GPSDD, IIASA, UN Women, UNCBD and UNCTAD to produce a Crowd4SDG policy brief based on recommendations that would also include good practices and success stories if identified during the study.



3. Methodology

This section describes the methodology used for this study.

A number of research questions have been formulated to guide the research:

- What do NSOs, the official statistics community as well as decision-makers at national and local levels know about citizen science in general and national level actors from Universities, NGOs, etc. working on citizen science or with citizen scientists? How does the official statistics community and decision-makers perceive the potential of citizen science data for SDGs?
- What are some examples of data gaps that citizen science could help address at the national and local levels from the standpoint of National Statistical Offices? Would citizen science data be more to some particular areas more than others, and, if yes, why? In what areas specifically can citizen science data be relevant for producing more timely or more granular data?
- What are some of the thematic areas where official statistics community and governments already attempt to use citizen science data? What are successful experimental projects that use citizen science data for SDGs or other country-level monitoring? Which indicators are governments most interested in exploring the potential of citizen science data for?
- What are the potential issues with data quality in the view of NSOs/the official statistics community and how can they be addressed (upfront where possible or in the aftermath if possible)?
- What are other obstacles in addition to quality that may impede NSOs/governments in leveraging citizen science data for monitoring SDGs?
- Why citizen science communities are thriving in some places and not others?

This study has used primarily a **qualitative research methodology** by examining the experiences of key stakeholder groups such as representatives of the official statistics community at national but also regional and international levels, and data users such as decision-makers from national and local governments and representatives of the citizen science community. It has also explored the emerging practices in this area, and involved a limited quantitative analysis at the level of perceptions using 5-1 or 3-1 scales for a selected number of aspects.

Data collection process involved two methods. First, <u>an online survey questionnaire</u> was sent out with 3 tracks for each of the key stakeholder groups: official statistics community, policymakers, and the citizen science community. The online survey was conducted between mid-February and mid-March 2021.

• <u>Track A</u>: sent to all Chief Statisticians and other members of NSO/NSSs (National Statistical Systems) as may be recommended by the Advisory Group, to Head Statisticians in International Organizations and Regional Organizations so they could



participate themselves but also invite relevant personnel to fill it in, and advertised at relevant events.

- <u>Track B</u>: sent to all SDG coordinating agencies inviting them to share with policymakers and planners in other Ministries as may be relevant, and to Local Authorities through UCLGs.
- <u>Track C</u>: advertised through citizen science community networks inviting representatives of citizen science communities to fill it in.

Key informant interviews were organized with the representatives of NSOs and other official statistics representatives who have been involved in citizen science or CGD projects, as well as those who would express their interest to be interviewed through the survey. The interviews started in June 2020 and continued until end March 2021, with most interviews with NSOs taking place after they have completed the online questionnaire.

Although various geographical regions were represented in the study – both in survey responses and in interviews, the survey featured a high number of respondents from some selected countries and no responses from others. This however was less of an issue given that the main survey objective was to collect qualitative data and the rankings were secondary. When repeated, qualitative data were reported only once. Another study limitation was related to the definitional issue. Citizen science data were understood sometimes too broadly to include digital trace and traditional household surveys. An analysis of qualitative responses has helped to identify such cases and exclude them from reporting. Finally, the response rate to Track C was quite low (12 respondents). As a result, these data were not used in the present report. Given that the main objective of the study was to understand the perceptions of NSOs and policy-makers and this third track was not part of the initial survey design, this decision has not affected the main objectives of the study and the study results. However, it is planned to undertake a separate analysis to understand how NSOs could engage with citizen science communities in the subsequent years of the project.

The study has benefitted from advice by the **study-specific Advisory Sub-Group** composed of:

- Gabriel Gamez, Luis Gonzalez Morales, Vibeke O Nielsen, Yongyi Min, Haoyi Chen, UNSD
- Karen Bett, GPSDD
- Linda See and Dilek Fraisl, IIASA
- Jillian Campbell, UNCBD (also member of Advisory Board of Crowd4SDG)
- Rosy Mondardini, CSCZ
- Sara Duerto Valero, UN Women
- Steve MacFeely, UNCTAD



4. Official Statistics: principles and quality assurance

The production and dissemination of Official Statistics by National Statistical Offices are guided by a number of principles and nationally specified quality assurance procedures or codes of practice.

The UN Fundamental Principles of Official Statistics (UNFPOS)⁶ adopted by UN Statistical Commission represent an overarching framework that is supposed to set the level of ambition and the standards to which all producers of official statistics – those from National Statistical Offices but also National Statistical Systems – have to aspire too.

The 10 UNFPOS include:

- 1. Relevance, impartiality and equal access;
- 2. Professional standards, scientific principles, and professional ethics;
- 3. Accountability and transparency;
- 4. Prevention of misuse;
- 5. Sources of official statistics;
- 6. Confidentiality;
- 7. Legislation;
- 8. National coordination;
- 9. Use of international standards;
- 10. International cooperation.

The activities of NSOs and the NSSs are in most cases regulated by national statistical legislation that states those principles, defines the roles and responsibilities within NSS, may specify its governance mechanism and its engagement with other stakeholders: data users, data holders, and more recently in some cases alternative data producers.

The quality assurance frameworks for official statistics serve to operationalize these principles. The quality assurance for official statistics is ensured through Codes of Practice or National Quality Assurance Frameworks that define quality assurance criteria. A lot of criteria are common across NSSs around the world and are also reflected in the UN Quality Assurance Frameworks for Official Statistics intended for UN agencies and in the UN National Quality Assurance Framework developed as a model for National Statistical Systems. The UN NQAF 2019⁷ is structured around: management of the statistical system, institutional environment, processes and outputs.

Managing statistical system:

- Coordination of the NSS;
- Relationships with data users, data providers and other stakeholders;
- Managing statistical standards.

Institutional environment:

- Professional independence;
- Impartiality and objectivity;
- Transparency;
- Statistical confidentiality and data security;

⁶ General Assembly Resolution 68/261. UN Fundamental Principles of Official Statistics. 29 January 2014.

⁷ UNQAF Manual for Official Statistics, Statistics Division, UNDESA, 2019.



- Commitment to quality;
- Adequacy of resources.

Processes:

- Sound methodology;
- Cost effectiveness;
- Appropriate statistical procedures;
- Non-excessive burdens on respondents.

Outputs:

- Relevance;
- Accuracy and reliability;
- Timeliness and punctuality;
- Accessibility and clarity;
- Coherence and comparability;
- Metadata.

Table 1 below shows the correspondence between UNFPOS and UNNQAF criteria. The columns contain UNFPOS from 1 to 10 mentioned above and the rows contain QAF criteria. The stars show which principle supports most strongly the related criteria in the row (mostly 1 per criteria), the circles show other supporting principles.

These principles and the quality assurance mechanisms have informed the design of this study and were used as a grid to analyze the results and develop recommendations.

In addition to the principles used by the official statistics community, one of the principles supported by the Office of the High Commissioner for Human Rights as part of the human rights-based approach to data is the principle of self-identification. This principle implies that the parameters of the population should be defined by the members of the population themselves and communicated via their (individual) decisions to disclose, or not disclose, their personal identity characteristics (e.g., their indigenous status, religion, or sexual orientation). In other words, the categories of identity should be developed through a participatory approach to enable an optimal engagement with data collection. Other principles that can be more directly mapped to the UNFPOS include participation, data disaggregation, transparency, privacy, and accountability⁸.

⁸ A Human Rights-Based Approach to Data. Leaving No One Behind in the 2030 Agenda for Sustainable Development. OHCHR Guidance Note to Data Collestion and Disaggregation. United Nations, 2018.



Quality principles		Fundamental Principles of Official Statistics								
		2	3	4	5	6	7	8	9	10
Level A: Managing the statistical system										
1: Coordinating the national statistical system								*		
2: Managing relationships with data users, data providers and other stakeholders					*			0		0
3: Managing statistical standards									*	
Level B: Managir	ng the	instit	utiona	l envi	onme	nt				
4: Assuring professional independence	0	*					0			
5: Assuring impartiality and objectivity		0	0	0	0		0			
6: Assuring transparency			*				0			
7: Assuring statistical confidentiality and data security						*				
8: Assuring commitment to quality		*								
9: Assuring adequacy of resources										
Level C: Managing statistical processes										
10: Assuring methodological soundness		٠			0				0	0
11: Assuring cost-effectiveness					*				0	
12: Assuring appropriate statistical procedures		*			0					
13: Managing the respondent burden					*					
Level D: Ma	anagi	ng sta	tistical	outpu	its					
14: Assuring relevance	*		0		0					
15: Assuring accuracy and reliability	*				0					
16: Assuring timeliness and punctuality					0					
17: Assuring accessibility and clarity			0							
18: Assuring coherence and comparability			0						0	
19: Managing metadata			*						0	
egend:										
 Fundamental Principles of Official Statistics (usually one) providing very strong support 										
Additional supporting Fundamental Principles of Official Statistics										

Table 1. UNNQAF principles and supporting UNFPOS Source: UNQAF Manual for Official Statistics, Statistics Division, UNDESA, 2019.



5. Main survey findings

Responses have been provided to questions under **three tracks**: official statistics community track, policy-makers track and citizen science community track. This section summarizes the key characteristics of the survey participants and main results for each of the three tracks.

5.1 Track A. Official Statistics Community

144 representatives of the official statistics community participated in the survey through the end⁹, most of them with 10 or more years of experience in this area (Graph 1). Different geographical regions were represented except North America (Graph 2).



Most of them had been members of National Statistical Offices. Some respondents also came from National Statistical Systems – other Ministries producing official statistics and local level – as well as from international and regional organizations (Graph 3). Around 20% of respondents had direct experience with or worked on non-traditional data sources¹⁰ as the main area of expertise, and another 20% had some indirect experience (Graph 4). Around 55% of those with some experience replied that they dealt with citizen science data (Graph 5).

⁹ 251 persons started answering the survey but this number dropped to 144 just after a few questions.

¹⁰ Administrative data and Earth Observation data were deliberately excluded from this list given their widespread use and that they were not the main focus of this study.

















Around 64% of participants were familiar with successful examples of data generated with support from citizens (Graph 6).

Among examples of citizen science data with which respondents were either familiar or their NSOs were engaged in included data generated from beach clean-ups, self-assessment for social data protection, adequate housing, access to food, gender rights, land cover monitoring project with volunteers work as ground-proofer, and Data bases collected from citizen to compute indicator 16.10.1 "Number of verified cases of killing, kidnapping, enforced disappearance, arbitrary detention and torture of journalists, associated media personnel, trade unionists and human rights advocates in the previous 12 months". In addition, crowdsourcing more generally and google maps have been mentioned too.



Close to 17% were aware of a project run by their Organization involving citizen data (Graph 7). Among factors that have sparkled such collaborations were collaboration with the University, dedicated capacity development projects in the field, awareness, the opportunities offered by such collaboration for timely data and the legal base. Based on the regional



breakdown for respondents having indicated the country, the use of citizen science data was particularly popular among respondents from the Asia and the Pacific region (Graph 8).



When provided with a potential list of **specific impediments** to the use of citizen science data, the **15 respondents who had experience with such data projects** ranked as highest:

- Limited access to data (67%);
- Legal issues with access or use of data (60%);
- Incoherent use or lack of use of statistical standard concepts, definitions and classifications undermining accuracy, reliability, coherence and comparability of the resulting statistics (53%);
- Selection bias due to data not being representative (53%);
- A lack of information about how the data are being produced (53%).

In qualitative feedback, the effective legal framework and institutional arrangements were identified as important preconditions for leveraging citizen science data and noted that those were often missing. The rules on access and confidentiality are typically those defined in general in the Statistical Acts, and some indicated that **their current legal basis has not provided a mandate for engaging with other stakeholders.** Some pointed to the need for the modernization of their statistical legislation. The strength of institutional coordination typically varies across countries with some NSO having a clear and well recognized mandate and effectively leading the NSS, and others – less so. This seems to also influence the opportunities and level and quality of engagement of the NSO with other stakeholders beyond NSS. Lack of cooperation with local authorities was specifically highlighted in the survey as an area that needs to be addressed to better leverage citizen science data from local level. Another limiting factor identified in the survey was the limited ICT integration/modernization in the process of data production and dissemination.

Some of the common solutions included brokering partnerships abroad to address capability constraints, introducing validation procedures for the data coming from citizens, conducting technical meetings and stakeholder workshops, using citizen science data in combination with other sources of information, and NSOs actively providing guidance to other institutions.



A question on main impediments to the broader use of citizen science data was asked to all survey respondents, incl. those who had no experience with data from citizens. 110 persons ranked them as follows:

- Lack of awareness (68%);
- Inability to ensure the use of statistical standard concepts, definitions and classifications (58%);
- Lack of methodological guidance (56%);
- Lack of human capacities to run experimental statistics projects (55%);
- Technological limitations (51%);
- Lack of technical and financial support to use citizen science data for the production of official statistics (50%);
- Statistical legislation does not enable/ prevents engagement (49%);
- Sustainability of access to the data source (49%);
- Fear of misinterpretation and misuse for citizen science data and other non-traditional data sources (47%);
- Fear of bias due to design of data collection (45%);
- Absence of vibrant citizen science community (41%);
- Normative loophole: no mechanism or protocols on key issues related to confidentiality, mandate, impartiality, etc. (40%);
- Irrelevance of citizen science data to main data requirements (40%).

In the qualitative feedback, in terms of barriers, one of the respondents noted that "sustainability of access is perhaps the main perceived barrier, making the upfront investment in this method of data collection something that NSOs are reluctant to take on at this stage". Other barriers included:

- The lack of clarity on what is possible, and how such data can make a contribution to filling in data gaps;
- In some cases, NSOs lack either technical capacity to lead on such projects or time to apply/initiate them;
- Sometimes, there may be resistance to introducing solutions where there is no clarity on how the quality issues can be addressed;
- Difficult dialogue with NGOs in one case as the latter do not consider statistics as part of their organizational knowledge due to the association with "the ideologically oriented statistics in the past (i.e., legacy of doubts about professional independence);
- Other problems included the discrepancy between the concepts and methodology used by civil society organizations;
- Missing legislation;
- Metadata are usually not available;
- Collection of citizen science data is not widespread enough, limited to some concerned people;
- Citizen science data /research cover / is carried out mostly for smaller territories or cities and rarely for the whole country.





The UN Quality Assurance Framework for Official Statistics identifies requirements along 3 areas: Enabling Environment, Outputs, and Processes. It also recommends that the national quality assurance framework be applied to all data and statistics produced outside of the national statistical system that are disseminated with the help and support of a member of the national statistical system or that are used for government decision-making.

In terms of the potential to address gaps on specific SDG indicators, respondents have highlighted the potential of citizen science data/CGD for all Goals with the ones with the highest ranking presented below:

- Goal 5. Gender equality;
- Goal 6. Water, sanitation and sustainable water management;
- Goal 13. Climate change;
- Goal 1. Poverty eradication;
- Goal 15. Terrestrial ecosystems;



- Goal 11. Sustainable cities and human settlements;
- Goal 3. Health;
- Goal 16. Peace, justice and institutions.

Where respondents felt such data source could be particularly useful were for:

- Measuring the access to and quality of public services;
- Informing environment (land / marine) indicators;
- Measuring poverty, in particular non-monetary, reflecting the locally rooted notion of poverty;
- Measuring what is typically not covered in the official statistics (the example provided was related to sex workers economy but may be relevant for other informal economy or yet illicit activities);
- Helping fill in gaps in official statistics that NSO/NSS cannot address due to budget constraints;
- Providing more timely data;
- Providing more granular data and data on vulnerable population groups;
- Adding complementary statistics;
- Reporting on some of the Goal 16 indicators;
- Providing more qualitative data required by some indicators, incl. data on subjective perceptions;
- Collecting data on certain phenomena where local knowledge is important, for example, knowing species.

Some respondents have identified specific indicators where citizen science data could contribute (Box 2).

Key findings can be summarized as follows. More than 20% of respondents had some experience with citizen science data although there is a lack of clear definitional borders on what exactly falls within citizen science data.

The key impediments as perceived by respondents with citizen science data experience compared to those without are somewhat different. Experienced respondents highlighted among key issues limited access to data, legal issues, incoherent use of statistical standard concepts, definitions and classifications, selection bias, and lack of metadata. This is in contrast to overall response showing the need for greater awareness, methodological guidance, human resource and technology limitations.



Box 2. Indicators mentioned specifically by survey respondents as those where citizen science data could contribute:

6.3.1 Proportion of domestic and industrial wastewater flows safely treated;

7.2.1 Renewable energy share in the total final energy consumption;

11.3.1 Ratio of land consumption rate to population growth rate;

11.2.1 Proportion of population that has convenient access to public transport, by sex, age and persons with disabilities;

14.3.1 Average marine acidity (pH) measured at agreed suite of representative sampling stations;

9.1.1 Proportion of the rural population who live within 2 km of an all-season road;

14.1.1 (a) Index of coastal eutrophication; and (b) plastic debris density;

16.1.2 Conflict-related deaths per 100,000 population, by sex, age and cause;

16.4.1 Total value of inward and outward illicit financial flows (in current United States dollars);

16.5.1 Proportion of persons who had at least one contact with a public official and who paid a bribe to a public official, or were asked for a bribe by those public officials, during the previous 12 months.

In terms of potential, official statics see opportunities for Goal 5 and gender, as well as environmental data, data on poverty, human settlements, and Goal 16.

5.2 Track B. Policy-makers

Track B was intended for policy-makers comprising a diverse audience from national governments, civil society, universities and academia, local governments, international organizations, and some private sector representatives. 126 respondents started answering the questionnaire, but the number dropped to around 84 respondents.

More than 60% of respondents were women and around 5% chose not to identify the sex. Close to 80% were national government representatives (Graph 10). The latter came from different ministries, in a number of cases from Ministries of Economy and Finance but also Prime Minister's Office, Ministries of Agriculture, Forestry and Water, Environmental Protection Agencies or regulatory bodies, Ministries of Labour, Social Protection, Home Affairs, Education, Minority Rights, Youth.





34.5% of respondents had more than 10 years of experience, another 29% and 27% situated themselves in the groups with 1 to 3 years of experience and 4 to 7 respectively (Graph 11).



Respondents have described a variety of tasks as data users for modelling and policy analysis but also for developing policies and reviewing budget execution, developing and monitoring progress on National Development Plans, SDGs and sectoral indicators, etc. More than 50% used data for policy-making or planning purposes, close to 40% used data for modelling, policy analysis or evaluation, and 42% were involved in programme management / had to communicate data as part of their responsibilities (Graph 12).



Among policy-makers, all regions were represented. Latin American and the Caribbean region was again most represented similar to Track A (Graph 13).





Close to 18% of respondents had **direct experience with non-traditional data sources** or had this as their main area of responsibility. Another 19% had some indirect experience.

The areas where such data were used in respondents' knowledge included: youth employment, citizen's perceptions regarding urban planning and development, real time information trends in some areas of sustainable development, to inform investment decisions based on SDG gaps, oceanographic and meteorological data in situ and that from network applications, satellite data, social media data, migration studies by NGOs and Universities, wastewater research to study the spread of COVID-19, data on air quality generated by citizens, biodiversity data gathering, GIS, scanner data and mobile phone data.

Several respondents had a very precise idea about citizen science data, a few confused it with social media data or with data produced for citizens.

Out of 69 respondents who answered the question, 19 were aware of **the currently ongoing projects in which their Organization was involved**. Those included projects related to data on refugees, on safe drinking water access, youth employment, air quality data, water assessments in dry regions, participatory poverty assessments, road death index. Approximately one third of respondents confirmed that their NSO was involved. In other cases, either project had validation mechanisms, or a scientific institution acted as a partner. Among factors explaining strong citizen engagement respondents named the facts that some problems directly impact citizens' lives, desire to strengthen Government's accountability, citizens concerned by high road death rates, awareness of the need for local action, and citizen's feeling ownership.

Respondents acknowledged that some projects had quality issues related to bias, coverage, quality, reliability or conceptual/methodological issues. The respondent having referred to bias in the context of water assessments in dry regions noted that there have been improvements in the quality of data with the time as a result of training. The quality of data from the project on air quality in schools that faced challenges with coverage and reliability had been addressed through triangulation.



According to this group, the top three factors preventing the broader use of citizen science data were the data quality considerations, the lack of awareness about opportunities it offers, and technological limitations in the country (Graph 14).



Among other impeding factors specified were legal regulations for access to information produced by the different levels of society with their due restrictions on confidential data, Incidence of political issues in collaboration and results (supposedly referred to issues with impartiality), lack of strong institutional arrangements, lack of statistical and institutional coordination, missing compatibility with existing/official data systems. One respondent noted possible limitations of citizen science potential to help collect data on vulnerable population groups that lack access to technology. Other noted the fact that citizen science data work well in communities with developed civil society which is a pre-requisite pointing to the lack of citizen participation in data collection culture as a major impediment.

In terms of potential, respondents felt that citizen science data could in particular be **useful** in providing more granular data or data on SDG indicators with gaps (Graph 15).





The perception of the potential of citizen science data for monitoring the SDGs was slightly different from that of NSOs. Goal 4 was ranked as first, Goals 1 and 13 as second followed by Goal 6 in the third place.



In terms of main findings, one can highlight that 37% policy side respondents having participated in the survey indicated to have had direct or indirect experience with citizen science data although some seemed to include plain social media data. Close to 28% indicated they were aware of the currently on-going projects in which their Organization was involved but only one third of them involved NSO. Interestingly, the lack of approval of data

D5.1 - Initial report on relevance and quality-related considerations of citizen-science generated data



by NSO was identified as one of the impediments (number 4). The top 3 impediments to the use of citizen science data identified by policy-makers were data quality considerations, lack of awareness and technological limitations.

In contrast to official statistics, policy-makers saw special potential for informing Goal 4 (Education) indicators in addition to those SDGs already identified by statisticians. They have also indicated that citizen science data could be particularly useful in providing more granular data and informing SDGs with data gaps.



6. Examples of citizen science data for SDGs and lessons learnt

The study has revealed several examples of the use of citizen science data for monitoring SDG indicators. Stakeholders have then been interviewed to discuss their experience in more detail and document some of the key findings and lessons learnt.

6.1 Key findings from the interviews

Among the interviewed NSOs, most had projects involving citizen science data. Institutionally, such projects were hosted by different Departments. Several developed countries and countries from Latin American countries run these projects under the experimental statistics portfolio either in combination with other data sources such as social media or new types of household surveys (see case study from ISTAT below). More recently, many started creating data science centers or departments. Such data are published as experimental statistics data which are not considered official statistics by some NSOs (unless they mature and move into regular production) and are considered official statistics by others. Others have such projects hosted directly under SDG monitoring units. In fact, SDGs data demand has become a significant factor that has spurred an interest in leveraging new data sources, including data from citizens, and in particular in countries that launched open SDG data portals. In such cases, data typically go through a quality assurance procedure by NSO and are then published as non-official data on the SDG portals. In other cases, such projects were hosted either in relevant sector statistics departments or in departments responsible for international cooperation where projects were funded by development cooperation agencies.

For some countries, the **lack of a legal basis** for NSO to engage with non-traditional data sources presents a major impediment. Most of the interviewed countries had the mandate to do so and felt it was an important enabling factor. In Colombia, for example, a recent decree has explicitly called on expanding the collaboration with alternative data sources.

Many have identified access to data as a key issue. In case of purely citizen science data, the challenge is often to know what kind of data exist and how to access it. As the pilot conducted by the Philippine Statistical Authority and PARIS21 described below shows, the need to conduct regular inventories and build up trust and confidence between citizen data holders and producers of official statistics (NSOs). Some countries indicated challenges related to accessing data from local governments. In case of mixed sources when citizens participate in the analysis of data from social media, access to social media data may be unequal. Some social media providers are changing their policies on access to API and make the service provision payable which may negatively affect such projects on financial or public procurement regulation grounds. In other experimental statistics projects involving mobile data and back transactions data, data are typically provided for free, but many telecommunication companies and banks are unwilling to share it or to share full datasets. One of the concerns for telecommunication companies is revealing the state of their operations to competing companies. **Strong partnerships** and the leveraging of data user engagement mechanisms here seemed to be of high relevance.

Related to access is the issue of **sustainability** of access to such data source. This also implies a defined frequency of production in the future. For many projects to move from experimental statistics to regular production, sustainability is critical. However, often this is beyond the control of NSOs when such data are produced with external funding – CSOs' funding depending on donors or that of international partners. Some experimental data do not need to become regular, however, but can still be useful at the moment when they are



produced such as, for example, the COVID-19 vulnerability index developed in 2020 in Colombia that served its purpose but would now need to be redefined.

The **relevance** of data didn't seem to be an issue overall. Several NSOs are exploring citizen science data specifically for addressing some of the data gaps on SDG indicators. A few others have pointed out however that they saw great potential in citizen science data when it comes to producing new indicators to provide some additional insights on the state and development of the country's society, economy, and environment.

When it comes to **confidentiality/privacy issues**, the views were mixed. On the NSO side, there are clear rules and well-established anonymization procedures. When data are being used/ published on NSO website, these conditions are met. It is not clear however whether confidentiality and privacy are respected during the data collection process by civil society / citizen science organizations unless this is clearly stated in the metadata. Whether this should be a concern or not for NSO is not sufficiently discussed.

A key requirement for NSOs to be able to use datasets produced by other stakeholders is a proper **metadata** fully describing the data set, how it was collected and what statistical procedures were applied. **Coverage** was identified as an important issue by many. Several interviewees noted that citizen science data often only concentrates on a small area, not even county, only villages, which makes its use hard. The projects designed ahead with stewardship of NSOs could address this issue at the design stage. In other cases, however, leveraging the existing initiatives may be the best way to go forward from the cost-effectiveness and timeliness perspectives or adjustments could be made to those with support of NSO. Several **methodological issues** may arise here. Overall, one of the main comments related to the use of new data sources was that NSO has to deal with data where outputs were not designed ahead and, therefore, the estimation techniques are very different. While this can hardly be circumvented in the case of Big Data and new estimation techniques have to be used here, this can be addressed in case of some of the citizen data projects if designed in collaboration with NSO or following a procedure recommended by NSO.

Several NSOs were working either on a **quality assurance mechanism** for data not initially collected by producers of official statistics and re-used for official statistical purposes or **guidelines** for civil society organizations and others involved in the production of citizen data. Several international organizations have developed or currently working on recommendations for quality assurance for citizen generated data. A major obstacle to leveraging citizen science data/CGD was a lack of know-how and **methodological guidance for NSOs** on how to develop and run such projects. The examples are presented in the case studies below.

The **perceived opportunities** related to citizen science data /CGD are multiple. Some NSOs focus on exploring how to <u>fill in data gaps on some SDG indicators</u>. Developed countries see significant potential when it comes to <u>environmental indicators</u>. Citizen science data can be a useful source to complement satellite imagery data and enable the measurement not only of the size but also of the healthiness of ecosystems. Many developing countries see potential related to <u>gender-based violence data</u> or data on <u>some "sensitive" indicators under Goal 16</u>. The Philippine Statistical Authority pilot has pointed to the potential of CGD for more than 80 SDG indicators in the country context. In particular, this is relevant for strengthening the availability of <u>more granular data</u>, as well as for countries that have <u>limited capacities to cover SDG indicators</u> through traditional surveys. Other NSOs pointed to the potential of citizen science data / CGD to bring <u>some new insights in different sector</u>



statistics, reduce costs and improve the timeliness of the data production¹¹ or fill in gaps that cannot be covered through traditional surveys such as informal economy, for example. Significant potential of citizen generated data for informing more qualitative indicators related more broadly to national sustainable development is acknowledged by many countries. Finally, a number of successful experiences of citizen engagement during humanitarian context points to a significant potential of citizens' contribution to data collection in situations of natural disasters or humanitarian crisis more broadly. Several pilots had been run by universities and international organizations involving citizens in ground-proofing during damage assessment but NSOs are not sufficiently familiar with this work. A number of the opportunities are perceived when citizen science is used in combination with other non-traditional data sources such as social media (see example from Mexico)¹². The examples are presented in the case studies below.

The interviews have also demonstrated that NSOs can use both "passive" and "active" approaches to leveraging citizen data. By passive here we understand approaches that are limited to NSO playing the role of standard-setter – through the development of guidelines for citizen science projects or CSOs for example - and scoping from time to time if data of sufficient quality becomes available/produced to inform selected indicators. By active here we understand approaches where NSOs proactively develop and manage projects and partnerships with citizen science initiatives, academia, CSOs and other stakeholders to inform relevant sustainable development – SDG or other relevant - indicators. The examples from UK and Kenya below would be closer to what we call here "passive" approaches, and the examples from Ghana, from Colombia on Goal 16, or from Mexico are the examples of "active" approaches.

6.2 Citizen generated data to monitor selected SDGs

The study has revealed several examples of the use of citizen science and citizen generated data for monitoring SDG indicators.

Case study 1. Environmental SDG indicators

Four environmental SDG indicators already benefit from citizen science data at the global level. Those include:

- Indicator 14.1.1: (a) Index of coastal eutrophication; and (b) plastic debris density;
- Indicator 15.1.2: Proportion of important sites for terrestrial and freshwater biodiversity that are covered by protected areas, by ecosystem type;
- Indicator 15.4.1: Coverage by protected areas of important sites for mountain biodiversity:
- Indicator 15.5.1: Red List Index.

A major benefit perceived by environmental statisticians is that citizen science data can help measure what satellite imagery cannot: the quality or healthiness of local ecosystems at a significant scale, beyond the capacity of many NSOs/National Environmental Ministries/Protection Agencies.

Source: Fraisl, et al. (2020), and interview with UNEP/UNCBD expert.

¹¹ An example of successful experimental projects where machine learning proved to be more accurate than human coding is related to the automated classification of economic activities and occupation. In Mexico, this has matured from an experimental project to regular production.

¹² Crowd4SDG project has also generated a dataset using this methodology to measure compliance of the public with mask-wearing and social distancing measures.



Case study 2. UK and ocean litter

UK Office for National Statistics (ONS) is looking into the use of the citizen science data from a charity (Marine Conservation Society) for compiling the indicator on ocean litter pollution: SDG indicator 14.1.1 part (b) on plastic debris density. The data source is currently undergoing the quality evaluation under the protocol developed by UK ONS SDG data sourcing team, discussed later in the report before it can be added on the SDG monitoring platform. This collaboration has been prompted following research by the SDG data team on possible sources from non-official data. In addition to the protocol approach and the open SDG reporting platform, another important factor of success is the sustained and targeted work of the SDG data team to identify possible data sources with advice from experts across ONS who can advise on new data sources. More information on source: https://www.mcsuk.org/how-you-can-help/

Source: Interview with UK ONS.

Case study 3. Gender-Based Violence in Ghana

Ghana's NSO, the Ghana Statistical Service (GSS), has developed a pilot project with the support of GIZ to build a data collection instrument – a new mobile application – enabling citizens to directly provide data on gender-based violence for 3 pilot districts. The project was a collaboration of Ghana's Statistical Service, Ministry of Gender, Children and Social Protection, and The Ministry of Local Government and Rural Development. The application, named 'Let's Talk Ghana' was developed using the design thinking approach through workshops with the involvement of several local IT companies who participated in the final bidding process as well as Civil Society Organisations, National and Local Government staff, and citizens themselves.

The questions were designed by the GSS in collaboration with the Ministry of Gender, Children and Social Protection to speak to four SDG indicators (5.2.1; 5.2.2; 11.7.2; 16.2.3) and the application provided an opportunity to report both personally or on behalf of another person. Confidentiality was ensured during the application installation process and the survey responses were submitted with no personal information. This was explained during the launch of the workshops in November 2020. The data received during the seven-week pilot data collection phase in November – December 2020 were then analyzed, with the records received before and on the day of the launch of the application excluded as test data. The GSS is in the process of finalizing the project report and will be comparing the results to the administrative data from the police records. Current police records are expected to experience underreporting from survivors who do not wish to initiate legal action.

The expectation is that the data collected through the application will better estimate the true prevalence of gender-based violence as barriers to reporting experiences are removed. One important aspect for the choice of the districts for the pilot was mobile network coverage, with the selected districts having varied access. However, a basic phone alternative was developed where users could dial a short code and receive a free callback giving the same questions in the local dialect. This aimed to combat biases from digital inclusion, mobile internet connection, and literacy concerns in public adoption of the solution. The GSS is moving the pilot into a wider roll-out phase, however, here some of the key factors will be: sustainability especially in the publicity required to promote technology solutions; data quality, and accessibility i.e. mobile coverage. The report will be published in 2021 by Ghana Statistical Service and hosted on their website.

Source: Interview with Ghana Statistical Service.



Case study 4. Freedom of speech in Colombia

A list of SDG indicators was defined between the Colombian statistics office DANE, OHCHR Colombia, and the Presidential Council of Human Rights. Based on this list, a diagnosis matrix was developed and last updated by the Presidential Council of Human Rights on 13/03/2021 to identify data sources for each relevant variable for the Goal 16 indicators prioritized. These include data from CSOs, administrative data from Government, and surveys.

The version 2.0 of the NSS – National Statistical System in Colombia - already includes by law the possibility to widen the scope of official data producers beyond the public sector, since the requirements to produce official data are based on the fulfillment of technical criteria, and not in the type of organization. Thus, given the availability of the conceptual and legal framework to work on the production of official statistics with different stakeholders, DANE has started to explore possibilities to work with different sources of information, including civil society.

For the compilation of the indicator 16.10.1, DANE has thought about the inclusion of data from human rights institutions of civil society. DANE have first contacted some NGOs such as Free Press Foundation – FLIP to explain the need for the indicator and, with the help from OHCHR, have begun a preliminary contact to Viva La Ciudadanía and the Alliance, a group of several NGOs to review possible ways to work jointly. In leveraging CGD for this indicator, a crucial step was to create a space for dialogue and bring around the table relevant Government entities and civil society organizations to analyze all data currently available from a technical perspective. One of the challenges DANE has identified is the difficulty to compare data from Government (administrative) sources and CSO data, however, it is aimed that all data available will be considered.

The NSO would like to leverage the ECP (acronym in Spanish for the Political Culture Survey) more as a primary data source to SDG indicator 16.b.1 and has developed a technical note explaining the need of the ECP survey since this survey doesn't have secured resources. This technical note was submitted for consideration of various entities to have funding. Also, for 16.b.1 DANE is running a project to get data from social media and to use natural language classification for having a measure that could also be compared with the obtained measure from the survey (if the ECP is funded). Here it's relevant to consider that the possibility of having the two measures from different methods that would allow DANE to calculate the bias of the social media data and also to contrast both exercises so the methodology with non-traditional sources could be improved and used in next years, ensuring the measure of the SDG indicator 16.b.1 regardless the funding of the survey.

A technical challenge in using CGD for Goal 16 indicators is related to uncertainty about the fulfillment of all quality criteria, since in most cases for new sources/methods, standards haven't been defined yet. The idea is currently to involve a third party - OHCHR – and learn from experiences of other countries, also to consider the work already conducted by UNSD that is consolidated in the citizen data toolkit, an input planned to be incorporated in these projects. Overall, DANE is aiming to "spread the word in the data community" so advice, ideas could also be considered.

Source: Interview with DANE and DANE documentation.



6.3 Citizen-generated data to provide key statistics relevant for the SDGs

One of the things revealed as a result of the study was that for many developing countries citizen generated data can be helpful in covering areas of economy and population segments that cannot be covered in their traditional business surveys or household surveys by design, more specifically on the sectors of the underground or informal economy.

Case study 5. Philippine pilot on CGD for official SDG reporting

The Philippine Statistic Authority has conducted a pilot study in 2019-2020 jointly with PARIS 21 to define the approach to leveraging citizen-generated data for official reporting on SDGs. The work comprised the send-out of a questionnaire, "Inventory of Data Holdings and Citizen-Generated Data of Civil Society Organizations (CSOs)/Non-Governmental Organizations (NGOs", to review and assess data collected by CSOs and how they are currently processed and used, in the context of the SDGs and the PDP. The results were mapped again the SDG data gaps identified previously by PSA. This was followed by the Launching Workshop that was an opportunity to develop technical guidance and reviewed quality assurance measures for possible use of CGD. This exercise has resulted in identifying 81 indicators to which CGD could contribute.

Under Goal 13, for example, for indicator 13.2.1. Number of countries that have communicated the establishment or operationalization of an integrated policy/strategy/plan which increases their ability to adapt to the adverse impacts of climate change, and foster climate resilience and low greenhouse gas emissions development in a manner that does not threaten food production (including a national adaptation plan, nationally determined contribution, national communication, biennial update report or other) the following CGD were available:

- Data on preparedness, adaptation, loss and damage, vulnerability, exposure, risk and resilience (including socio economic profile of respondents).

Source: PARIS 21, PSA & PSRTI Report "Use of CGD for SDG reporting in the Philippines: A case study", June 2020.

Case study 6. Kenya and sex workers economy

Kenyan NSO used civil society data for measuring the sex workers' economy. This is an area not covered in the household or business surveys. To enable the NSO get the data from them, the best approach was to use Civil Society Organization so that they can connect them to the sex workers. These enabled the NSO to administer their research questions with ease. The coverage was ample, as it covered major hot areas where sex workers are found. The coverage may not entail all sex workers since for those who are not registered with CSO and may be hard to approach and find them. Confidentiality was adhered as questions asked were not revealed to the CSO and only the NSO used the data for statistical purpose.

It is expected that the use of such data will benefit from the promulgation of the new guidelines for CSOs to help improve the quality of citizen generated data and make it more apt for use by the NSO. As a way forward, the NSO looks forward in working with various CSOs to fill in the data gaps such as the underground economy.

Source: Interview with Kenyan NSO.



6.4 Citizen-generated data and local level / social monitoring of the SDGs

The potential of CGD for local level monitoring and social accountability mechanism is well acknowledged¹³. Main challenge for NSOs who produce data about the country is that such data have limited coverage. Even when such data produced in many regions, if concepts and definitions applied in these different regions are not consistent between themselves, it is impossible to aggregate such data. They can however serve useful purposes for decision-making at local level or national policies that involve geographically targeted interventions or resource allocation in particular when comparisons between regions become possible or when the region had already been selected for a policy intervention. In addition, as the below example from Ghana shows, local level monitoring can help strengthen social accountability providing a direct feedback mechanism for service providers and increasing the overall transparency around access to quality basic services.

Case study 7. Colombia and social cartography

Colombia's NSO DANE is providing advice for a resilience-enhancing project called UKRI GCRF Understanding Risks & Building Enhanced Capabilities in Latin American Cities with a focus on social and participatory cartography, led by the University of Warwick (UK) together with partner universities and institutions in Colombia and Brazil. The objective is to co-develop a framework for risk and vulnerability monitoring in informal settlements to build a disaster risk monitoring platform, for residents and government agencies to collaboratively implement. The framework design involves interviews with citizens and feedback from academia and NSO on possible indicators to measure social and physical vulnerability. This approach involves citizens not only in the data collection but also in defining what vulnerability is and how to measure it.

Importantly, this approach involved the citizens not only for data collection but also aimed to empower them by co-productively defining the vulnerability and how to measure it. Initial results point to differentiated concepts and methodologies for disaster risk monitoring, and the utility of alternative and participatory approaches for generating riskrelated data, jointly validated by communities and government agencies.

This trans-disciplinary pilot project to which DANE contributed with expert advice focused on Medellín and Rio de Janeiro and comprised geographers, statisticians, data experts, geologists, community activists, and engineers. After project completion, the aim is to further develop this methodology for enhanced disaster risk monitoring and ultimately scale up to other municipalities across Colombia and Latin America in general with histories of internal displacement.

Source: Interview with DANE and DANE documentation.

¹³ See, for example, the "Basic principles of Community-Based Monitoring" developed by United Cities and Local Governments and The International Observatory on Participatory Democracy in 2014. <u>https://www.uclg.org/en/media/news/basic-principles-community-based-monitoring</u>



Case study 8. Ghana and local waste management

Inefficiencies in solid waste management pose a significant environmental challenge for Metropolitan, Municipal, and District Assemblies (MMDAs). However, the essential statistics needed for making informed decisions and policies are lacking in the country. On this premise Ghana's NSO, the Ghana Statistical Service, also ran another innovative pilot project combining data collection with a social accountability/feedback mechanism. <u>CleanApp Ghana</u>. A GPS-based application was developed through a design thinking process to enable citizens in three pilot districts to report on waste mismanagement such as uncollected household bins or overflowing community containers. The reports included a location, photo, and estimated volume and were received by both the Local Government's Environmental Health Officers and the private solid waste service provider responsible for collection, who also had an area on the app to mark issues as resolved.

The by-product data from the app and basic phone short code versions were accessed by the NSO via a database with the aim to compute SDG indicator 11.6.1. The pilot faced many challenges in incentivizing the service providers to participate; technical problems and estimation errors given by citizen reporters. Despite this in the seven-week data collection period of the pilot, over 900 valid voluntary reports were made through CleanApp Ghana.

Source: Interview with Ghana NSO.



6.5 Citizen-generated data beyond the SDGs

Producing data for SDG indicators is important of course but citizen science data can be useful for measuring society's and planet's well-being and the strength of the economy in new ways, increasing the timeliness of data availability, reducing costs and the respondent burden.



Case study 9. Mexico sentiment analysis

Mexican NSO, INEGI, has set out to develop a strategy to fill in the gaps where traditional data fail, especially with SDGs. It has set up a new data scientists' team in October 2019 looking at social media, bank transactions, satellite imagery, phone data to make some indicators that are hard to measure.

One of the experimental projects they have been running since 2014 on the sentiment analysis involved the use of social media data and volunteers' help for classification. The social media component was managed by INEGI while citizen science collaboration was brokered through a partnership with the Universidad Tec Milenio. The sentiment data are available every day at the national level and state levels. They are published as experimental statistics which are not considered official statistics by INEGI.

Students are asked to label manually tweets provided by INEGI as negative or positive. Each tweet is labelled by more than 10 people to ensure accuracy. During the first-year pilot, there was an issue with some tweets being labelled wrongly but the cross-referencing by 10 persons helps identify malicious behavior by rare individuals and counter it. Tweets are then cleaned and normalized and the manual classification is used to help build training datasets with more than 50, 000 values. The machine learning algorithm developed with the support of data scientists from two research centers, INFOTEC and CentroGeo, is then used to label millions of tweets automatically. The georeferencing allows producing sentiment analysis index not only at the national level but also at the level of states. Correlations are clearly observed for some events such as the 2017 earthquake with more of negative tweets on average while Christmas appears to be the happiest day of the year.





Case study 9. Mexico sentiment analysis (continued)



The key advantage offered by the partner University was also the fact that it has campuses across the country. This has allowed capturing differences in local slang when labeling the tweets.

Students were motivated to participate in this project for several reasons: doing social media analysis was perceived as an exciting task, collaborating in a brand new, innovative project with INEGI which was perceived as very prestigious, and the specialization of the University in technology and engineers preparing students with strong IT skills across different educational programs.

One of the reasons for the success of the project was easy access to the Twitter data. INEGI used to collect all possible data from the country in geolocalized tweets. This is however about to change as Twitter is modifying their access to API policy restricting the number of inquiries in 1 day. The essential service for the project will therefore become payable and would imply the application of highly regulated public procurement procedures.

Source: Interview with INEGI and <u>https://www.inegi.org.mx/app/animotuitero/#/app/multiline</u>.



Case study 10. Citizens contributing to Italy's trusted smart surveys

As part of its experimental statistics portfolio and a broader European Statistical System initiative to set up an EU smart surveys platform, Italy is rolling out trusted smart survey pilots combining traditional survey sampling techniques and sensor data from mobile devices. Italy is taking the lead on the initial design effort for the development of an EU platform as part of the European Statistical System (ESS) project.

This builds on the work carried out on the applications and online data collection as part of the Time Use Survey and Household Budget Survey and can inform on these topics as well as Living Conditions, Health.

Some of the key issues that the pilots by participating EU countries aim to address include effective tactics to engage and involve "citizens"; and accuracy and comparability of (trusted) smart survey data relative to regular data collection. Citizens will be engaged through incentive schemes, in particular pilots will look at respondent recruitment and motivation strategies. Another important point pilots aim to elucidate is generic and respondent concerns about privacy.

Smart surveys are considered mixed-mode surveys integrating data collected from different sources. A broad range of tools is considered, including apps, wearables, a sensor of mobile devices. Machine learning is used to deal with huge volumes of sensor data.

The results of the project are expected to be presented by mid-2022 in a final dissemination event.

Source: Interview with ISTAT.



7. Guidelines for citizen data producers

A number of countries have been working on preparing guidelines - and in some cases formal Quality Assurance or validation mechanisms - for non-official data sources to encourage these alternative data producers to improve the quality and meet some minimum standards set by NSOs.

Case study 11. UK Protocol for non-official data for monitoring the SDGs

UK's Office for National Statistics has an open SDG National Reporting Platform with most of the data coming from official data sources. To address gaps on some of the indicators, the Office for National Statistics SDG team was open to leveraging new data sources and decided to put in place a protocol for non-official data specifically for monitoring the SDGs.

The protocol is work in progress. The draft consists of an ethical gateway and a scoring matrix. The Ethical Gateway serves as a pass/fail mechanism with three criteria all of which have to be met:

- Ethics and Privacy
- Transparency and Accountability
- Need

Once this condition is met, a scoring matrix is applied to evaluate the given data set using the following criteria:

- Relevance (0-3)
- Methods (0-3)
- Coverage (0-3)
- Timeliness (0-3)
- Data journey awareness (0-3)
- Quality assurance (0-3)

The average score is then calculated and 1.5 points are used as a threshold to decide on whether or not the source could be used for the SDG platform.

The non-official data protocol was developed based on and is aligned with the UK Statistics Authority <u>Code of Practice</u> and its voluntary application procedure applied to non-official data sources. The latter are then listed on UKSA portal as the Organizations voluntarily applying the Code. The 3 key pillars of the Code are Trustworthiness, Quality and Value.

In addition to data, research articles can also be linked to using a similar protocol with some adjustments. The accessibility is one additional "must" criteria included on the ethical gateway for research articles, and the scoring matrix includes an additional dimension of publication quality. 1.5 average score is used again as a threshold but only the top 3 articles are referenced.

Source: Interview with and documentation provided by UK ONS



Case study 12. Colombia and Quality Assurance guidelines for experimental statistics

Colombian NSO, DANE, does not have a quality assurance mechanism for non-official data sources per se but has developed quality assurance guidelines for experimental statistics by adapting its quality assurance framework for official statistics. The experimental statistics workstream was launched in 2020 as a flagship initiative guided by the dedicated Technical Committee chaired by DANE's Director and Chief Statistician in Colombia and composed of all technical directors and some advisors.

At the moment all data created by DANE are official statistics. Experimental statistics are considered official statistics in Colombia according to Decree 2404 from 2019. They offer new ways of quantifying phenomena relevant for sustainable development and can be helpful in ensuring disaggregated information, address data gaps, combine traditional data sources such as censuses and surveys with new data sources, develop in-house capacities to run experimental projects. Experimental data need to meet a defined set of quality assurance criteria that was developed based on UN QAF for official statistics with some criteria dropped, as they didn't fit the specificity of experimental data.

Those criteria include:

- Relevance: extent to which statistics meet real users' needs;
- Accessibility: how easily statistics can be accessed by various users;
- Interpretability: how easily users can analyze statistics / how clear it is;
- Transparency: information accompanying the access to data by users (data + metadata);
- **Coherence:** extent to which the used concepts, applied methods, and produced results are logically connected;
- **Timeliness:** time span between the point in time during which the phenomenon is quantified and publication of statistics to ensure that it can be useful for informing decisions.

Compared to the UN QAF, an important missing criterion is the one related to accuracy and reliability. This distinguishes experimental statistics from regularly produced official statistics. The former often requires improvements on standards, coverage, and methods, and has not reached a sufficient degree of maturity to get into regular official statistics production processes; so reliability without defined standards can't be measured.

Before the introduction of quality assurance for experimental statistics, it has been hard to recognize and share this work. DANE had several innovative projects but the results couldn't be published in the absence of a quality assurance approach. The 2019 revision to the statistical law has also provided the mandate for the Colombian NSS to engage and incorporate alternative data sources in their official data production processes.

Source: Interview with DANE and DANE website



In addition, several international organizations have produced guidelines specifically on the use of citizen-generated data. This section highlights 3 of them: the PARIS 21 recent publication for NSOs on CGD, the GPSDD guidelines for Government and other organizations on CGD and the work of UN Statistics Division on CGD.

Case study 13. PARIS 21

Drawing on the experience of the pilot with the Philippine Statistical Authority for which PARIS 21 provided support, the Partnership has developed a set of recommendations for NSOs. The five core recommendations include:

- Communicate a working definition of CGD
- Identify the purpose for using CGD
- Implement quality standards for CGD
- Develop institutional capacities to coordinate with CGD producers
- Establish data repositories to facilitate data sharing

The proposed quality assurance approach is similar to UK ONS's approach that combines a set of quality criteria, the scoring matrix, and threshold level.

The quality criteria proposed by PARIS 21 for CGD include:

Figure 3: Quality a	ssessment framework for evaluating CGD
()	RELEVANCE: Degree to which CGD serve to address a specific purpose sought by the NSO
đ	ACCURACY: Measurement precision of CGD to estimate the actual values of the targeted phenomena (coverage, sampling method, non-response, etc.)
Ŕ	CREDIBILITY: Lewel of trust and objectivity of CGD (no attachment of interest pushed by the data producer)
J	TIMELINESS & FREQUENCY: Attribute of CGD to be up-to-date and available on periodic basis
\square	ACCESSIBILITY: Ease with which CGD are presented, released and made available to users
go	INTERPRETABILITY: Simplicity with which users can understand and use CGD correctly
- 60	COHERENCE: Comparability over time and across geographical units

Source: Reusing citizen-generated data for official reporting. A quality framework for national statistical office-civil society organisation engagement, PARIS 21, February 2021.



Case study 14. GPSDD and a multi-target audience guide on how to choose and engage with CGD

Together with Open Knowledge International and Public Data Lab, GPSDD has developed a guide "Choosing and engaging with CGD" intended for governments, international organizations and other stakeholders whish to run CGD project and therefore is not limited to NSOs.

The three key aspects examined in this Guide include the considerations related to:

- whether CGD is fit for your purpose;
- what type of participation/ contributions by citizen and what depth are expected;
- what are the larger workflows within which your data generation takes place.

The guide is structured along following questions:

- 1. What are your objectives, questions and data needs?
- 2. How can the engagement and participation of people help?
- 3. What resources are available to support CGD?
- 4. How can CGD be made public?
- 5. What considerations are relevant for data protection?

The guide contains a lot of useful and practical examples, and illustrates the various roles that citizens can play with regards to CGD:





Case study 15. UNSD

UN Statistics Division is currently developing a guidance to help countries better use of citizen-generated data for public policy. More specifically, the guidance includes three stages of work:

- 1. Developing a quality assurance toolkit for producers of citizen generated data to plan and document the process of data collection, processing, analysis and dissemination. When properly done, the documentation helps CGD producers better communicate the quality of collected data with the official statistical community, hence increases the likelihood of the CGD data being used to inform policy actions.
- 2. Developing a Toolkit for NSOs to incorporate citizen generated data in the official statistics.
- 3. Developing a total error framework for citizen-generated data by working with survey statisticians to estimate the total error of citizen generated data and to improve the accuracy of the overall measurement by identifying ways to correct the biases.

Source: Interview with UNSD.



8. Overall recommendations

Given the diversity of approaches to how citizens can contribute to data production but also how NSOs can engage with it and for what purposes they can use it, a flexible but comprehensive framework may be required.

In active approach projects, NSOs can provide data stewardship from the outset and have more control over the data production process: follow UN Fundamental Principles of Official Statistics, apply GSBPM, and ensure statistical outputs meet quality requirements. In purely citizen data projects – not involving Big Data – outputs can be designed ahead or at least shaped and the standard statistical concepts and techniques can be applied resolving a lot of methodological issues that arise in passive approach projects and in Big Data experimental projects.

In passive approach projects, the quality of statistical outputs can be controlled to see if they meet minimum requirements and to some extent statistical processes depending on what is available in the metadata. Here, some potential issues may arise on products, process and environment side: those covered in quality criteria applied for outputs by many NSOs and mentioned in the PARIS 21 guidelines and the UN Fundamental Principles of Official Statistics, such as confidentiality and impartiality.

Many NSOs may be interested in leveraging both approaches. Therefore, a comprehensive but flexible framework could be applied.

8.1 Strengthening capacities of NSOs to leverage citizen science data

Based on the lessons learnt and to boost the capacity of NSOs to leverage both active and passive approaches, the following areas need to be addressed:

- Updating if necessary the legal basis to ensure NSOs have the right mandate to engage with CSOs, Academia and communities;
- Strengthening partnerships with CSOs, Academia and communities who may potentially contribute to data production. This may include those who already produce data (passive approach) and those with whom NSO could engage on collaborative projects. Working with data user community can help identify not only those Organizations that may already be producing data but also those who may be interested in collaborating on new projects. For projects involving social media analysis with citizens' help, signing MoUs with social media providers could be something to explore;
- Defining clearly for what citizen science generated data will be used as recommended by PARIS 21. This will also be relevant for collaborative projects. GPSDD diagram showing the different roles citizens can play would be extremely useful in this sense;
- **Defining quality criteria** based on the many examples that emerged (discussed more in next section);
- Introducing a protocol / quality assurance mechanism. This can be a scoring matrix with threshold values as in the UK ONS example;
- **Providing training and capacity development for stakeholders** involved in citizen data production to enhance statistical literacy (for CSOs), improve the knowledge of the principles of official statistics (also important for Academia) and the awareness about the needs and the work of NSOs.



Box 3. UN National Quality Framework for Official Statistics (2019)

Managing statistical system:

- Coordination of the NSS
- Relationships with data users, data providers and other stakeholders
- Managing statistical standards

Institutional environment:

- Professional independence
- Impartiality and objectivity
- Transparency
- Statistical confidentiality and data security
- Commitment to quality
- Adequacy of resources

Processes:

- Sound methodology
- Cost effectiveness
- Appropriate statistical procedures
- Non-excessive burdens on respondents

Outputs:

- Relevance
- Accuracy and reliability
- Timeliness and punctuality
- Accessibility and clarity
- Coherence and comparability
- Metadata

Source: UNNQAF Manual, 2019.

For international organizations:

- Develop **methodological guidance** on how to run collaborative projects based on pilots and leveraging the expertise of NSOs, Universities and citizen science community.
- Promote peer learning between countries on CGD
- **Support projects in this area** that also help build local / NSO capacity to replicate / run on their own in the future.

8.2 Possible criteria to be included in guidelines / protocols for CGD

The main criteria used by a number of countries (UK, Colombia) with some variation and close to those proposed in PARIS 21 guide include:

Accessibility – (Anonimized) datasets should be easily accessible online to the broader public. Depending on a country context, having data published in local languages may be an added advantage.

Timeliness, Frequency and Sustainability -- Data should be available/disseminated on time to be used as evidence for decision-making (e.g., in humanitarian context, the rapidity of



access to data is critical). An added consideration would be the frequency of data production for those indicators where trends over time are important.

Accuracy and Reliability – Data should be produced using sound statistical procedures and methods to ensure the closeness of estimates to true values and the previously estimated values when preliminary figures are disseminated.

Coverage – Incomplete coverage may be a serious obstacle to an effective use of data. When no proper sampling techniques are applied, it may be difficult to establish to what extent incomplete data coverage are representative of the population or a given group. For national indicators, complete coverage for the country's population or territory is often a must. In some cases, incomplete coverage may be addressed through calibration or integration with other sources such as household surveys. Where it is not possible to correct bias, footnotes and metadata must be added to the estimates to warn the user to interpret the results with caution.

Relevance – It is important that the collected data are relevant to decision-makers to make progress on national development objectives and / or inform specific public policies. For NSOs using or wanting to use citizen data for SDG monitoring, nationally relevant/adapted SDG indicators are an important benchmark.

Metadata – The provision of a proper metadata and its accessibility together with the dataset is a key condition. Without metadata, NSOs cannot judge the dataset's compliance with many other criteria such as coverage, relevance, accuracy and reliability, coherence, comparability, and integrability. The access to metadata is also essential for data users so that the data can be used effectively and appropriately.

Coherence, Comparability and Integrability – The coherent use of standard statistical concepts and methods enables NSOs to combine datasets in different ways and make comparison across regions, over time. They allow data aggregation and its use in combination with other data sources. In some cases, proxy indicators can prove helpful when data for the indicator itself is not available, however NSOs can promote the coherence and comparability of the produced data through guidelines, publishing indicators with their metadata and training for CSOs.

Four additional criteria could be considered:

Documented data collection/production/dissemination process - Metadata should provide full information not only on the data set but also on the process of data collection. This is an important indicator in UK's protocol called data journey awareness.

Impartiality - The Organization supplying the data should be able to demonstrate that it is committed to impartiality in the data production process. While this may be less of an issue for Universities and scientific community, it may be less obvious for NGOs that combine advocacy, service provision and monitoring mandates. It is important they are committed to training personnel involved in data collection and production on statistical procedures, techniques and principles. The demonstration of impartiality can be addressed through proper metadata showing the application of sound statistical procedures and fully transparency on the process and outputs, but it is important to ensure metadata reflect the reality.

Confidentiality/Privacy – For datasets on citizens, the Organization supplying the dataset should demonstrate that the data collection process has involved a full consent from the



respondents during the data collection process / no violation of data protection legislation if such exists. As some Organizations may be collecting this data as part of their service delivery, such data may not be anonymized in their internal databases similar to administrative data. However, it is important to ensure that the respondents are fully aware and consent to data collection to avoid unethical behavior and that the published datasets comply with confidentiality/privacy.

Self-identification – An important principle to consider from the Human Rights based approach to data about individuals is the principle of self-identification allowing the person to define his/her gender, ethnic, cultural and other identities in accordance with his/her per perception and if he/she decides to disclose his/her identity. Population parameters should be defined by the population representatives through a participatory process.





9. Comparative analysis of data quality approaches between official statistics and academia/citizen science and recommendations

In this section, we compare the quality assurance framework used by the official statistics community as described in detail in previous sections and the one used by the academia.

The commonly used data quality criteria identified through a dedicated study Cichy C. and Rass S. (2019)¹⁴ in technology and scientific circles are cited below:

- **Completeness:** the extent to which data are of sufficient breadth, depth and scope for the task at hand.
- Accuracy: the extent to which data are correct, reliable and certified.
- **Timeliness:** the extent to which the age of the data is appropriate for the task at hand.
- **Consistency:** the extent to which data are presented in the same format and compatible with previous data.
- Accessibility: the extent to which information is avail-able, or easily and quickly retrievable.

In addition, a number of additional criteria emerge from the citizen science community are defined in another study Anhalt-Depies C.¹⁵ such as

- **Data quality**: this should address the quality of citizen science data, incl. accuracy and coverage, and address lack of standardized sampling protocol, poor spatial or temporal re-presentation, and insufficient sample size
- **Privacy**: to manage risks related to data collection encroaching on data privacy (e.g., images of identifiable persons)
- **Resource security**: prevent the possibility of data collection creating a risk for the security of natural resource (e.g., wildlife through location information)
- **Transparency**: scholarly research should be open to the public given the public resources sued to support it.
- **Trust**: building trust between scientists and the public is essential to citizen science.

Some of the criteria that seem to be less present or at least not explicitly articulated in the scientific data quality approach and which may need to be strengthened in citizen science projects aiming to contribute to monitoring SDGs and NSO data production include:

Relevance:

This is one of the areas which distinguishes significantly official statistics from academic data research. Official statistics requires that the data produced should be of direct relevance to data users – policy-makers and the needs should in fact be defined in dialogue with the latter. While there has been more recently more emphasis on policy-science interface overall and citizen science projects have been traditional looking at issues that are of growing policy relevance – environmental data and statistics – this criterion needs to be acknowledged more clearly for projects aiming to produce SDG relevant data.

Frequency and sustainability

Another key distinction of official statistics is that it aims to ensure for most of the indicators a specified frequency of production leading to time series data. This may be a challenge for

¹⁴ Cichy C. and Rass S., An Overview of Data Quality Frameworks, IEEE Access, 2019.

¹⁵ Anhalt-Depies, C., Tradeoffs and tools for data quality, privacy, transparency, and trust incitizen science. Elsevier, 2019.



some of the citizen science projects and is a key impediment to NSOs engaging more with citizen science data where the sustainability of access cannot be ensured. However, this may be less of problem for data on punctual or temporary issues, for example, in the context of disaster related damage assessment. Timeliness is key to making sure that the data can be used by policy-makers. In fact, one of the reasons why NSOs and policy-makers may be interested in leveraging citizen science data is because it can be made available faster. However, it is important that data is not only produced in a timely manner but also disseminated in a timely manner to serve the purpose of public policy making.

Metadata:

A detailed description of the dataset is a key requirement in addition to making the key produced statistics available in open access to the public. While this is done for some of the citizen science data sets, this should be a routine procedure for all projects aiming to make data useful for monitoring the SDGs. Metadata should be published together with the dataset in an open access in a timely manner.

Coherence, Comparability and Integrability:

In official statistics, the coherent use of standard statistical concepts and methods is essential to allow for comparison between countries, country regions and in time. This is an important element to consider for citizen science data producers for the data that they hope to be useful for monitoring global or national SDG indicators. For local level monitoring this may be less relevant when the challenge is specific to the country or location. Further, as the world evolves and new challenges arise, some of the new concepts may evolve within scientific community prior to or in parallel to these issues being defined by official statistics community.

Documented data collection/production/dissemination process:

For datasets produced outside of official statistics, it is important that NSOs could evaluate the quality. This would require that all stages of the data process – collection, production and dissemination – are well documented.

Impartiality:

The data production process should be objective, free of any political influences, supported by the application of sound statistical procedures and fully transparency on the process and outputs.

Self-identification:

This principle that came from the human rights perspective should be applied to data about individuals who should have the right to self-define themselves (e.g., ethnicity) and be able to participate in the definition of the identification parameters.

When citizen science datasets are assessed on their potential use for the monitoring of the SDGs and data production by NSOs, it would be important to pay special attention to the above criteria that had been so far less explicit in the citizen science data work and would need to be strengthened. This is a key recommendation of this report for citizen science community wishing to support the monitoring if the SDGs and produce data more likely to be used for informing public policies.



10. Conclusions and outlook

The proposed criteria are being and will further be tested during the data usability assessments on datasets provided by Crowd4SDG partners and those that will be coming out of the GEAR cycles. The first assessment is presented in Crowd4SDG Deliverable 5.2 " Data usability assessment and recommendations for SDGs for GEAR cycle 1 ». Further refinements will be considered based on this experience and the evaluation of new datasets in subsequent years of the Crowd4SDG project.

The findings of the report and criteria have also been presented and discussed during a panel discussion organized on 27 May 2021 with the Advisory Group members and several NSOs. The recommendations and criteria will provide a basis for a policy brief.

In addition, a more detailed **methodological guidance** can be developed on how to run collaborative projects in the context of possible pilots with NSOs and other UN and non-UN partners leveraging the expertise of NSOs, Universities and citizen science community.



11. References

Articles, manuals and publications:

Anhalt-Depies, C., Tradeoffs and tools for data quality, privacy, transparency, and trust incitizen science. Elsevier, 2019. https://reader.elsevier.com/reader/sd/pii/S0006320719301958?token=0E3F8A988B060275 75A7FAC5012037E7C2FB2D3E248FA80B34097022EACBB25DF30A008A3005CEBE4E49E8 0D2E01CA7D&originRegion=eu-west-1&originCreation=20210615174328

Bonney R, Ballard H, Jordan R, McCallie E, Phillips T, Shirk J, Wilderman CC, Public Participation in Scientific Research: Defining the Field and Assessing Its Potential for Informal Science Education. A CAISE Inquiry Group Report. 2009. Available at: <u>https://eric.ed.gov/?id=ED519688</u>

Cichy C. and Rass S., An Overview of Data Quality Frameworks, IEEE Access, 2019. Available at: <u>https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=8642813</u>

Crowd4SDG Deliverable 5.2 "Data usability assessment and recommendations for SDGs for GEAR cycle 1", June 2021.

Fraisl, D., et al., Mapping citizen science contributions to the UN sustainable development goals, 2020. Available at: <u>https://link.springer.com/article/10.1007/s11625-020-00833-7</u>.

GPSDD, Choosing and engaging with Citizen-Generated Data. Guide, GPSDD, Open Knowledge International and Public Data Lab, 2018. Available at: <u>https://www.data4sdgs.org/resources/choosing-and-engaging-citizen-generated-data-guide</u>

Haklay M., et al, Contours of citizen science: a vignette study, Royal Society Open Science, 2020. Available at: <u>https://osf.io/preprints/socarxiv/6u2ky/</u>.

Irwin A., Citizen Science and Scientific Citizenship: same words, different meanings?, 2015

Negri V., Scuratti D., Agresti S., Rooein D., Scalia G., Fernandez Marquez J.L., Ravi Shankar A., Carman M. and Pernici B. (2021). Image-based Social Sensing: Combining AI and the Crowd to Mine Policy-Adherence Indicators from Twitter, ICSE - Track Software Engineering in Society, May 2021. Available at:

https://re.public.polimi.it/retrieve/handle/11311/1161146/584481/ICSE_SEIS_Image_based _Social_Sensing%20%2837%29.pdf

PARIS 21, Reusing citizen-generated data for official reporting. A quality framework for national statistical office-civil society organisation engagement, PARIS 21, February 2021. Available at: <u>https://paris21.org/sites/default/files/2021-02/CGD_FINAL_reduced.pdf</u>

PARIS 21, PSA & PSRTI, Use of CGD for SDG reporting in the Philippines: A case study. 2020. Available at: <u>https://paris21.org/sites/default/files/inline-files/PSA-report-FINAL.pdf</u>

UCLG and IOPD, Basic principles of Community-Based Monitoring, 2014. Available at: <u>https://www.uclg.org/en/media/news/basic-principles-community-based-monitoring</u>

UK Code of Practice. Available at: https://code.statisticsauthority.gov.uk/



General Assembly Resolution 68/261. UN Fundamental Principles of Official Statistics. 29 January 2014. Available at: <u>https://unstats.un.org/unsd/dnss/gp/FP-New-E.pdf</u>

UNSD. UNNQAF Manual 2019. Available at <u>https://unstats.un.org/unsd/methodology/dataquality/un-nqaf-manual/</u>

United Nations. A Human Rights-Based Approach to Data. Leaving No One Behind in the 2030 Agenda for Sustainable Development. OHCHR Guidance Note to Data Collestion and Disaggregation, 2018. Available at:

https://www.ohchr.org/Documents/Issues/HRIndicators/GuidanceNoteonApproachtoData.p df

Links:

DANE. Experimental Statistics. Available at: <u>https://www.dane.gov.co/index.php/estadisticas-por-tema/estadisticas-experimentales</u>

INEGI. Available at: https://www.inegi.org.mx/app/animotuitero/#/app/multiline

Ghana NSO. Let's talk Ghana. Application. Available at: <u>https://play.google.com/store/apps/details?id=net.aoholdings.letstalk</u>

Ghana NSO. Clean up Ghana. Application. Available at: <u>https://play.google.com/store/apps/details?id=com.cersgis.ahonedie&hl=en_US&gl=US</u>

UK ONS & Marine Conservation Society. Charity for compiling the indicator on ocean litter pollution: SDG indicator 14.1.1 part (b) on plastic debris density. Available at: <u>https://www.mcsuk.org/how-you-can-help/</u>



Annex: list of abbreviations

Abbreviation	Description
AI	Artificial Intelligence
СВІ	Challenge-based Innovation (in-person coaching)
CBIx	Challenge-based Innovation (remote location)
CCL	Citizen Cyberlab
CGD	Citizen generated data
CS	Citizen Science
CSSK	Citizen Science Solution Kit
10	International Organization
GEAR	Gather, Evaluate, Accelerate, Refine
NSO	National Statistical Office
NSS	National Statistical System
017	Open Seventeen Challenge (online coaching)
SDG	Sustainable Development Goal
UCLG	United Cities and Local Governments
UNFPOS	UN Fundamental Principles of Official Statistics
UNNQAF	UN National Quality Assurance Framework
UNQAF	UN Quality Assurance Framework