# Response to "Evaluation of Indoor Localisation Systems: Comments on the ISO/IEC 18305 Standard"

Nader Moayeri

*National Institute of Standards and Technology, Gaithersburg, MD, USA*

## Abstract

Potorti et. al. wrote a paper [1] that was dubbed as the first critical reading of the international standard ISO/IEC 18305 [2], Test and evaluation of localization and tracking systems. The author of this paper served as the Editor of ISO/IEC 18305. As such, I have unique insight into and perspective on how the standard was developed. The comments made in [1] fall into two categories from my viewpoint. The first category are the comments that are worthy of consideration in a potential revision of the standard. The second category are the comments that I do not agree with. This paper is a detailed, point-by-point response to the comments made about ISO/IEC 18305 in [1]. The paper also provides guidance on the logistics for a potential revision to the standard and how others can get involved in that process. The views presented in this paper are primarily those of the author.

## Keywords

localization and tracking, indoor localization, test & evaluation, standards

## 1. Introduction

Performance evaluation of Localization and Tracking Systems (LTSs) is a challenging measurement problem for many reasons. First, the assumptions under which various LTSs function can be very different from each other. For example, in designing and deploying an LTS for a shopping mall, one can afford to install a lot of electronic equipment in the mall and carry out extensive measurements to "train" the system so that it would yield the desired level of localization accuracy and consistency. On the other hand, an LTS intended to provide situational-awareness for the firefighters going into a burning building and the incident command post at the scene is quite different from the LTS for the shopping mall. In this case, one cannot assume there is any electronic infrastructure in the building to enable localization. Even if such equipment existed, it would most likely be unavailable, because the first thing firefighters do at a fireground is to cut off electricity and natural gas supply to the building. Therefore, firefighters can rely only on whatever equipment they take to the scene and quickly deploy. As another example, an LTS intended for asset tracking in a large warehouse is different from the two LTSs just discussed, because in this case the LTS is supposed to locate objects that do not move on their own and perhaps some delay in estimating the asset location is acceptable, if it makes the estimate more accurate. Therefore, one has to be very careful when comparing different LTSs. They cannot be

compared based on localization accuracy only. One has to take all sorts of other factors into account.

Second, one has to test an LTS in several large buildings. The building size matters, because the degradation in LTS performance does not manifest itself in a small building. To have an appreciation for this requirement, consider the drift of an inertial measurement unit that accumulates over time as one walks large distances or the challenges of Radio Frequency (RF) ranging over several walls. Other factors that do matter include construction material(s) and whether the building is subterranean or a high-rise. Larger buildings tend to use more metal in their construction that poses a challenge to any RF-based or magnetic localization technique. What one should not do is to claim that an LTS works well based on testing in just a single-family house made of lighter construction materials. It would be highly desirable if an LTS could be tested in a laboratory according to procedures that could be replicated anywhere around the world. Unfortunately, it is not possible to test an LTS in a laboratory. It has to be tested in the field, which means in several buildings in this case.

Third, whether the Entity to be Localized / Tracked (ELT) can move by itself or not and how it moves has an effect on the LTS performance. Fourth, an LTS may employ all sorts of sensors and techniques. Hence, it is important to test an LTS in scenarios that pose a challenge to these sensors and techniques. For example, any vision-based systems using optical cameras would suffer in environments with poor lighting. Therefore, any LTS that uses an optical camera should not be tested in only ideal lighting conditions.

ISO/IEC 18305 adopts a black-box testing approach. It assumes nothing is known about the inner workings of the LTS under test. Any approach that attempts to tailor the testing to how the LTS under test functions, even though highly desirable for the developers of that system, would not be scalable, because there are just too many systems out there. ISO/IEC 18305 is meant to be used for testing several LTSs in a given set of buildings according to specified testing scenarios, and then comparing the performance of the systems tested. Unfortunately, this may not say a whole lot about how these systems would perform in buildings other than those used in testing.

Given all the complexities described above, it would be natural for experts in indoor localization to not fully agree on every aspect of LTS testing. That's why I said at the outset that this is a challenging measurement problem. The purpose of this response is to make those gaps narrower, if not to bridge them entirely.

The rest of this paper is organized as follows. Section 2 goes over the genesis of ISO/IEC 18305. In Section 3, we address all the points raised about the standard in [1]. We use *Italic* font for quotes from [1], followed by my view on that specific issue in normal font. If a direct quote from [1] includes a citation, that citation is to the reference list in [1] and not the one in this paper. In Section 4, I describe the mechanisms through which experts who were not involved in the development of ISO/IEC 18305 may participate in a potential revision of the standard. Naturally, there has to be a need for substantive changes in order to initiate a revision process. Finally, Section 5 presents the concluding remarks for this paper.

## 2. The Genesis of ISO/IEC 18305

*The standard is one of the results of the dissemination and exploitation plan of the EU FP7 project EVARILOS – Evaluation of RF-based Indoor Localisation Solutions for the Future Internet [10].*

Response- The development of ISO/IEC 18305 started in ISO/IEC JTC 1/SC 31/WG 5 in September 2012. ISO/IEC JTC 1/SC 31/WG 5 stands for International Organization for Standardization / International Electrotechnical Commission Joint Technical Committee 1 / Subcommittee 31 / Working Group 5. JTC 1 focuses on Information Technology, SC 31 focuses on Automatic Identification and Data Capture, and Working Group 5 focused on Real Time Locating Systems. WG 5 no longer exists, and the standards under its purview are now handled by WG 4, which focuses on Radio Communications.

This effort was separate from and independent of Project EVARILOS, as it has also been stated in Section 2.5 of [3]. Some researchers from EVARILOS got involved in the development of ISO/IEC 18305 in 2014. The standard benefited from their technical contributions, particularly to Appendix B, but the extent of their involvement was minimal. The standard development process was not a sub-project of EVARILOS. The standard was published by ISO on November 1, 2016.

## 3. Main Response

*The standard does not consider other T&E and evaluation purposes, like those that would be useful to system developers or the testers themselves.*

Response- The standard has already a few performance metrics that are primarily of interest to system developers. Two examples of such metrics are the mean of error vector, which represents the systemic bias of the LTS, and the covariance matrix of the error vector. In addition, susceptibility and resilience, which are optional performance metrics for LTS use in mission critical applications, are more likely to be of use to system developers than to users.

It is not clear why testers need separate evaluation procedures.

*The standard gives a good overview on the architecture of common localisation systems, methods and sensors used, and an interesting appendix that does a good job of enumerating various failure modes of many sensors. Both are rather comprehensive, but for some reason lack any references to radio-based Device-Free Localisation (DFL) systems [13]. This omission is especially noticeable given that the standard is particularly attentive to fire-fighter scenarios, where DFL systems, especially tomographic ones [14], are in principle attractive.*

Response- It was a conscious decision by the developers of 18305 to restrict the scope to LTSs that rely on the ELT to carry or be equipped with a localization device. DFL does not belong to that category.

I regard RF-based DFL and tomographic imaging as *niche* areas. They are more appropriate for intrusion detection than for localization. Even if such a system detects the presence of a person at certain location, it would be hard for it to determine the identity of that person. Imagine the situation where it's needed to locate an unknown number of people in a large room based on changes in the RF signal propagation environment. That would be difficult with DFL of references [13] and [14] in [1]. Yet, thermal imaging cameras are used by firefighters to

detect presence of people behind a wall.

A more promising approach is DFL via a network of optical cameras with overlapping fields of view installed in a building, such as an airport terminal, for surveillance purposes. Such a system would rely on face recognition algorithms to recognize the person whose image is captured and then on tracking algorithms to track the person's movements. Then a handover would take place as the person leaves the field of view of one camera and enters the field of view of the next one. There is extensive R&D activity and benchmarking in that area. The evaluation of such vision-based systems needs a separate standard. It would have not been easy to handle that topic or other types of DFL in 18305.

*A standard-compliant system test should perform a separate test for each scenario. This is a heavy requirement for some systems. For example, Table I lists how many scenarios are defined for each of the 15 combinations of ELT type and building type.*

Response- References should have been made to Table II in 18305 and not to Table I. Table II provides *guidance* on which building types and which scenarios should be used for testing an LTS meant for a particular ELT type. The standard allows the possibility of using fewer scenarios for testing if requested by the system manufacturer, as long as it is made clear in the T&E report for that system which scenarios were used for testing.

Nevertheless, I acknowledge that the burden of testing an LTS according to 18305 is high. We struggled to balance the need for comprehensive testing, as explained in Section 1, against the requirement to keep LTS testing cost manageable. Maybe a second look at the required T&E scenarios is warranted in a potential revision to 18305.

*It is interesting to note that apartments are missing from the list of test buildings. Distinguishing an apartment in a building from a single-family house is significant, both because of building differences (materials, structure) and the interaction between similar systems in adjacent apartments that should be evaluated in the case of apartments.*

Response- I agree that indoor localization for apartments is an important issue, because many people live in apartments. For example, in responding to an E-911 call, it is important to know from which apartment in the building the emergency call originated.

Testing in apartment buildings requires some thinking. Interactions between similar systems has been mentioned above. There is also the possibility of having one system for a building housing many apartments. One should take into account that localization accuracy may vary from one apartment to another. For example, the system may perform better in an apartment on the side of the building as opposed to one deep inside. If testing is to be done throughout the entire building, then this may not be all that different from testing in an office building, which is already part of 18305.

*The metrics considered by the standard are various statistics applied to point errors, that is, to errors measured from the ELT position to a series of test points which is defined as the ground truth. The standard does not consider other approaches related to the overall trajectory, such as the Fréchet distance [15]. We think that the reason behind using point error statistics instead of comparing trajectories [16], [17] is that the latter is less adequate to navigation purposes, for which the real-time identification of the position is more important than the path followed.*

Response- We were not aware of the Fréchet distance. Hence, its use was never considered in the development of 18305.

One has to justify the use of Fréchet distance by explaining in which application(s) it would

arise and then decide whether it makes sense to add it to 18305. One needs to ask how often this distance is used in practice. I agree that this distance might be useful when a person or an autonomous robot traverses a path and it matters how far apart the actual path taken is from the LTS's estimate of that path.

*While it is true that computing two statistics or thirty requires essentially the same amount of work, we believe that the usefulness of a standard is greatly improved if only a small set of results is produced, so that it is clear what is important for characterising a given system.*

Response- Localization accuracy is an important performance metric for an LTS, but it is not the only one. Other important performance metrics include the size and weight of the ELTD (the device carried by the ELT), its battery life, location update rate, system capacity in terms of how many ELTs it can localize / track simultaneously, and whether the LTS needs the floor plans of a building in order to operate. These are discussed in Clause 11 of 18305. Given that there are so many dimensions to LTS performance, would it be objectionable to use more than one metric for accuracy?

Many people use the mean of magnitude of horizontal, vertical, and/or 3D error as the performance metric for accuracy. In the Positioning, Navigation, and Timing (PNT) community, they is more interest in Root Mean Square (RMS) values and x% points on various Cumulative Distribution Functions (CDFs).

In 18305, we presented all the performance metrics used by various communities to measure accuracy. In principle, I agree that it would be desirable to reduce the number of such metrics in 18305 in a potential revision of the standard, but I do not believe it is a good idea to artificially reduce that number to one. I doubt that a single performance metric for LTS accuracy would be appropriate for all LTS applications and would capture the entire picture.

*Some of the metrics to which we object are very useful for system debugging, others for system tuning after installation and periodically afterwards, others are appropriate to validate the correctness of the testing procedure itself: in none of these cases are the metric useful to the final user of the system.*

Response- 18305 was written with the LTS users in mind primarily, but not exclusively. As mentioned earlier, some of the performance metrics are more appropriate for system developers. One possibility is to take them out of 18305. Another way to proceed is to divide the metrics into different categories depending on whether they are more useful to users, system developers, etc. A combination of the two approaches is also possible.

At this point, it is not clear to me which performance metrics are appropriate for system debugging, tuning after installation and periodically afterwards, or to validate the correctness of the testing procedure itself.

*In order to improve the credibility of the test report, the standard should mandate that the presentation of numerical results should include not only the number of samples, as already specified, but also the amplitude of a 90% confidence interval for each figure. The amplitude of the interval shall be discussed, a reasonable requirement would be for the interval to be not wider than $\pm 20\%$ around the reported value.*

Response- This is a good suggestion, but its implementation needs to be specified and it has to be ensured that it makes sense mathematically. The problem is that when people compute confidence intervals, they use sampling with/without replacement. This makes sense for i.i.d. data. LTS error is spatially and temporally correlated.

*As a last note on metrics in general, the standard should recommend minimum and maximum performance requirements for various types of systems.*

Response- I do not agree with this suggestion at all. 18305 shows how to measure things. It stays away from setting benchmarks, which should be left to policymakers. One example of a policy is the requirement set by the Federal Communications Commission (FCC) in the US that requires wireless service providers achieve 50 m location accuracy or better for 80% of the E-911 calls by April 2021.

While the lifetime for a standard for LTS testing may be measured in decades, technology advances rapidly and with it LTS accuracy levels feasible improves.

*Probabilities of correct detection: These are simple probabilities computed on floors or zones, the latter in case zones of interest are defined. For some use cases, this may be all that is needed, for example in the common case where one only needs to know the position of the target at the room level.*

Response- This is a reasonable suggestion. However, a few issues need to be considered before one rushes to make changes in 18305:

1. Are there other LTS applications for which a simplified testing procedure is warranted?
2. It is still important to test in different building types. Therefore, the guidance on building types and size is still applicable. So is the guidance on testing scenarios.
3. Determining floor number is straightforward if the building looks like a layered cake with the same thickness for all layers. In other words, when all the floors in the building have the same shape as its footprint and all floors have the same height. However, there are many buildings that do not meet these requirements. Examples of such buildings have been given in 18305. For such buildings, one needs to have a computer program that unambiguously and consistently determines the floor number. This is hardly trivial in the general case. There are also complications about floor numbering, because a building may have half-floors or the number of floors may be different depending on from which side one looks at the building. These issues are not trivial.
4. Similar problems exist with determination of room number. With a "human in the loop", the problem is easy. One can just overly the horizontal location estimate provided by the LTS on the building floor plan and the human would tell you this is in Room *blah blah*. However, an LTS does not have a human in the loop. It would need a robust classifier that would map the horizontal location estimate to room number. This may require a decision tree, whose design would not be trivial.
5. These issues have already been mentioned in 18305.

*These statistics are mostly inappropriate as the result of a system-level, black-box approach to test and evaluation of a localisation system, for various reasons.*

Response- Perhaps a few of these metrics can be removed or placed in a special category in a revision of 18305. However, I disagree that these statistics are mostly inappropriate. How can one claim the means of magnitudes of horizontal and vertical errors are inappropriate? They are widely used in scholarly papers. In addition, I have already mentioned that the PNT community uses RMS values as performance metrics.

Once again, 18305 is *primarily* intended for users, but not *exclusively*.

*Similar reasoning can be applied to almost all of the above metrics. Apart from the root mean squares, the first- and second-order statistics above should not be presented.*

Response- A wholesale statement like this regarding a long list of performance metrics is not acceptable. I have already argued why some of the metrics are useful.

*An overall 3D bias which is significantly different from zero points to an installation problem. Indeed, as mentioned in the standard, it is normally null, and if different it indicates an installation or coordinates measurement problem. It may also indicate a system's weakness. In any case, while extremely useful to the system developer, to the installer and to the tester, this information is useless or even misleading from the point of view of the final user. It should not be presented.*

Response- I have already addressed this issue. The most logical approach is to put a few performance metrics in a separate category intended for system developers or installers.

*Root mean square of error magnitude may be an exception, because in fact it gives a good grasp on the interesting characteristics of the error with a single number. However, it is less useful than the quantiles of error, which are discussed next.*

Response- I have already mentioned that the PNT community uses both. Let's not summarily discard RMS values.

*First, it should be clear that a generic 3D error, in xyz coordinates measured using any reference system, while often easier to measure, is not generally useful.*

Response- That is too strong a statement to make. I disagree.

*When the ELT is an object, height may matter, but again floor is the most important information, so height should be relative to the floor.*

Response- Asset tracking in warehouses, where assets are typically stacked, is an important application of indoor localization. In that case, height is the only thing that matters, because warehouses are typically single-floor buildings.

Why rush to discard a metric when we don't yet have a standard for floor numbering let alone robust software to determine floor number?

*The standard prescribes the computation of floor detection probability, as mentioned in section II-C1: this is necessary but not sufficient. There are some ways to deal with this problem, which are summarised in table II. The standard should prescribe the use of the second and preferably also the third solution mentioned therein.*

Response- It has been acknowledged in Table II of [1] that the floor penalty in "2D Euclidean distance with floor penalty" is somewhat arbitrary. I agree. I doubt there will ever be consensus in using that approach. In a sense, it is like forming a figure of merit for an LTS that is a linear combination of scores for different aspects of an LTS, such as horizontal accuracy, vertical accuracy, battery life for localization device, and many other factors that have been described in Clause 11 of ISO/IEC 18305. Even if we agree that a linear combination is acceptable as opposed to combining these scores in a nonlinear fashion, there will never be agreement on the weight for the linear combination. The weights would vary from one LTS application to another.

The idea of using the Real Distance in Table II [1] is clever, but one has to be careful before adopting it as a measure of LTS accuracy. In order to compute the Real Distance, one has to run a shortest path algorithm in a building to find the shortest distance that a person has to travel (using corridors, stairs, and perhaps elevators) to get from the LTS estimate for the ELT location to the real ELT location. This metric tends to penalize buildings that have complicated floor plans. This is easiest to see in case of a floor plan that looks like a labyrinth, which would

make the Real Distance excessively large when, for example, the true ELT location is 2 m way from the ELT estimate, but a person would have to travel 50 m on a shortest path to get from the estimated location to the true location. Would that be fair?

It is worthwhile to discuss Real Distance further. If the community feels that it is superior to other measures of accuracy in certain cases, then it can be added to the standard with the necessary qualifications on when it can be used.

*Better yet would be to use more quantile values, so to give a good approximation of the CDF (cumulative distribution function). The standard should require to use four quantiles: 0.5, 0.75, 0.9, 0.95, which makes it easy to compare two systems and to set minimum requirements for normative reasons.*

Response- The standard already suggests that the CDF of an error magnitude should be plotted. If one has the CDF, then one can get from it not only the above four quantiles, but any others that may be of interest.

*The standard recommends that 50–100 test points are set up per floor, and that at least half of them is used for each scenario.*

Response- The standard specifies the density at which test points should be deployed in a building. Specifically, it states that "A test point shall be deployed in every 5-10 $m^2$ of area in the single-family house described in 10.1.2 and in every 50-100 $m^2$ of area in each of the other building types used for LTS T&E."

*This means that, if high quantiles are to be measured, the course should be walked through at least several times to meet the requirements mentioned in section II-C.*

Response- Repetition is an important factor in design of experiments. It is somewhat addressed in the use of location-specific accuracy as a performance metric, but no T&E scenarios have been included to test for location-specific accuracy. There were two considerations in this decision. It was felt that the burden of going through all the T&E scenarios was already high, and we did not want to make this worse by requiring repeating each scenario in each building type a few times. We also felt that by using a sufficient density of test points in each building, perhaps a spatial average of accuracy over the test points would play the same role as a time average implied by repeating the scenarios.

*Latency is of great importance for static object localisation, that is, cases when the ELT is not normally moving. Systems dedicated to this use case may adopt a long measurement procedure in order to filter out noise and trade promptness for accuracy, for example by giving few estimates per minute or even at longer time intervals.*

Response- This is not always true. For example, there is a requirement in the US that a location estimate must be available at a Public Safety Answering Point (PSAP) within 30 seconds of when an E-911 call is made. Therefore, latency may matter even in cases where the ELT is stationary.

I agree with the statement in the next paragraph in [1] that latency is not explicitly important in certain cases. For example, when a person needs navigation inside a building, an LTS typically updates the person's location at certain rate. If the LTS has large latency, the effect would be reduced accuracy, because even if the LTS generates very accurate location estimates, by the time it is made available to the user, he/she may have moved by several meters.

*Walking means human walking at a speed of about 5 km/h, which is in fact too fast for a generic indoor environment.*

Response- That is a fair point. We had read somewhere that the average human walking speed in 5 *km/h*. I suspect that speed is for outdoors. We need to find out what the average indoor walking speed is.

## 4. Mechanism to Revise ISO/IEC 18305

A standard developed by the International Organization for Standardization and the International Electrotechnical Commission (ISO/IEC) typically undergo a five-year review cycle to decide if it should be reaffirmed, revised, or withdrawn. As of May 2021, a fifth-year review for ISO/IEC 18305:2016 has not yet started but could begin soon. Also, a revision project can be initiated at any time. For those who may be interested in participating in a revision of ISO/IEC 18305:2016 and who have not participated before in the national body processes of ISO and IEC, a snapshot is provided below.
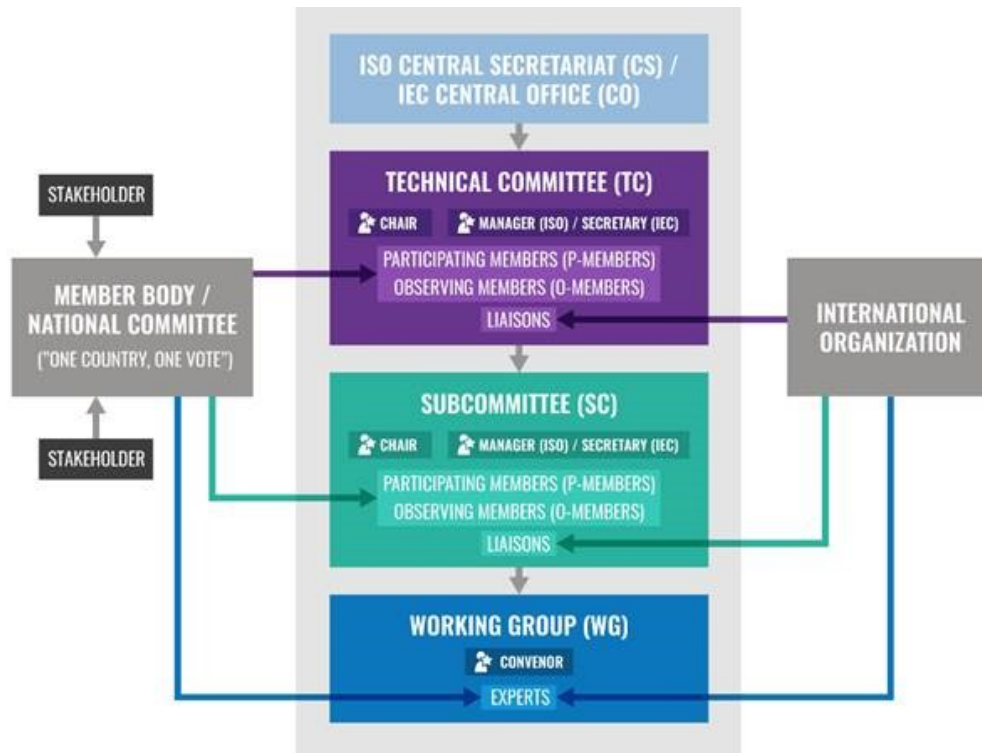
ISO and IEC are private sector international standards developing organizations. Each has only one member per country, typically a national standards organization. Individuals or companies can't become a member of ISO or IEC. They can only participate in ISO groups via their national mirror committee or in IEC groups via their national committee.

Joint Technical Committee 1 on Information Technology (JTC 1) is a joint committee of ISO and IEC. ISO/IEC 18305:2016 was developed and is maintained by ISO/IEC JTC 1/SC 31 Automatic identification and data capture techniques [4]. JTC 1/SC 31 is administered by ISO. Consequently, national body participation is through national mirror committees.

The role of the national mirror committees is to develop their national consensus positions on all technical and administrative matters brought before the respective ISO Technical Committee (TC) or Subcommittee (SC). ISO TCs and SCs usually form Working Groups (WGs) composed of experts to develop their standards. The national mirror committees nominate their national experts to participate in each WG. Figure 1 summarizes how one can participate in ISO/IEC standardization.

The JTC 1/SC 31 home page lists its standards, structure, members, and liaisons. Presently, there are 24 Participating Members and 26 Observing Members. There are also some Liaison Organizations, which must have an international membership to be considered for liaison. The member listings include contact information for each national standards organization. Contact your national standards organization to find out about your national mirror committee for JTC 1/SC 31.

Aside from the mechanics of revising an ISO/IEC standard described above, it would be highly desirable to have several experts involved in the deliberations so that the decisions are not made by just 1-2 WG experts. It would be great to have representations from all stakeholders, such as the industry, user community, and possibly academia and organizations dealing with policy and regulations. It would be even better to have multiple people from industry and user community involved, because companies focus on different LTS technologies and users should represent different LTS applications or at least the most important ones.

**Figure 1:** Participation avenues in ISO/IEC standards development

## 5. Concluding Remarks

This rebuttal presents the views of the Editor of the international standard ISO/IEC 18305 [2] on the issues raised in [1] regarding the standard. Responses have been provided for all the points raised in [1]. The Editor of 18305 agrees that some of the comments made in [1] are worth discussing in a potential revision to the standard. He also explains why he disagrees with the other comments.

The rebuttal also describes the process through which an ISO standard may be revised. Anyone can make the initial request for the standard to be revised by making the case that a revision is needed. More often than not, the Editor for the original version of the standard is involved in a revision process. One possible way to proceed for others interested in seeing the standard revised is to send an email to the Editor so that he can gauge the level of interest and the extent of modifications suggested. Individuals interested in seeing the standard revised should be prepared to participate in the revision process at the national level in their own countries and at the international level within ISO/IEC JTC 1/SC 31.

# References

[1] F. Potorti, A. Crivello, P. Barsocchi, and F. Palumbo, "Evaluation of indoor localisation systems: comments on the ISO/IEC 18305 standard," Nineth International Conference on Indoor Positioning and Indoor Navigation (IPIN), pp. 1-7, Sept. 2018.

[2] "ISO/IEC 18305:2016, Information technology – Real time locating systems, Test and evaluation of localization and tracking systems," Nov. 2016. https://www.iso.org/standard/62090.html

[3] T. Van Haute, E. De Poorter, I. Moerman, F. Lemić, V. Handziski, A. Behboodi, A. Wolisz, N. Wirström, T. Voigt, P. Crombez, and G. Glorioso, "D2.4 Final Version of the EVARILOS Benchmarking Handbook," 2015.

[4] "ISO/IEC JTC 1/ SC 31 Automatic identification and data capture techniques." https://www.iso.org/committee/45332.html