



TECHNICAL REPORT

**Speech and multimedia Transmission Quality (STQ);
Best practices for robust network QoS
benchmark testing and scoring**

Reference

DTR/STQ-00219m

Keywords

3G, benchmarking, data, GSM, LTE, network,
QoE, QoS, scoring, service, speech, video,
ViLTE, VoLTE

ETSI

650 Route des Lucioles
F-06921 Sophia Antipolis Cedex - FRANCE

Tel.: +33 4 92 94 42 00 Fax: +33 4 93 65 47 16

Siret N° 348 623 562 00017 - NAF 742 C
Association à but non lucratif enregistrée à la
Sous-Préfecture de Grasse (06) N° 7803/88

Important notice

The present document can be downloaded from:

<http://www.etsi.org/standards-search>

The present document may be made available in electronic versions and/or in print. The content of any electronic and/or print versions of the present document shall not be modified without the prior written authorization of ETSI. In case of any existing or perceived difference in contents between such versions and/or in print, the prevailing version of an ETSI deliverable is the one made publicly available in PDF format at www.etsi.org/deliver.

Users of the present document should be aware that the document may be subject to revision or change of status.

Information on the current status of this and other ETSI documents is available at

<https://portal.etsi.org/TB/ETSIDeliverableStatus.aspx>

If you find errors in the present document, please send your comment to one of the following services:

<https://portal.etsi.org/People/CommitteeSupportStaff.aspx>

Copyright Notification

No part may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying and microfilm except as authorized by written permission of ETSI.

The content of the PDF version shall not be modified without the written authorization of ETSI.

The copyright and the foregoing restriction extend to reproduction in all media.

© ETSI 2019.

All rights reserved.

DECT™, **PLUGTESTS™**, **UMTS™** and the ETSI logo are trademarks of ETSI registered for the benefit of its Members.

3GPP™ and **LTE™** are trademarks of ETSI registered for the benefit of its Members and of the 3GPP Organizational Partners.

oneM2M™ logo is a trademark of ETSI registered for the benefit of its Members and of the oneM2M Partners.

GSM® and the GSM logo are trademarks registered and owned by the GSM Association.

Contents

Intellectual Property Rights	6
Foreword.....	6
Modal verbs terminology.....	6
Introduction	6
1 Scope	7
2 References	7
2.1 Normative references	7
2.2 Informative references.....	7
3 Definition of terms, symbols and abbreviations.....	8
3.1 Terms.....	8
3.2 Symbols.....	8
3.3 Abbreviations	8
4 Governing Principles for Mobile Benchmarking	8
4.1 General	8
4.2 Fair Play	8
4.3 Comparing networks with different coverage extents	9
4.4 Comparing networks with differing technology use	9
4.5 Test device selection	9
4.6 Test server selection.....	9
4.7 Test method transparency.....	10
4.8 Advice and best practice for web-page selection	10
5 General Description.....	11
6 Test Areas.....	12
6.1 General	12
6.2 Geographical divisions	12
6.2.1 Cities.....	12
6.2.2 Roads	12
6.2.3 Complementary areas	13
7 User Profiles.....	13
8 Test Metrics.....	13
8.1 Introduction	13
8.2 Telephony.....	13
8.2.1 General.....	13
8.2.2 Telephony Success Ratio	13
8.2.3 Setup Time.....	14
8.2.4 Listening Quality	14
8.3 Video Testing	14
8.3.1 General.....	14
8.3.2 Video Streaming Service Success Ratio	14
8.3.3 Setup Time.....	14
8.3.4 Video Quality.....	14
8.4 Data Testing	14
8.4.1 General.....	14
8.4.2 Success Ratio	15
8.4.3 Throughput	15
8.5 Services Testing	15
8.5.1 General.....	15
8.5.2 Services.....	15
8.5.2.1 Browsing	15
8.5.2.2 Social Media	15
8.5.2.3 Messaging	15

8.5.3	Success Ratio	15
8.5.4	Timings	16
9	Weighting	16
9.1	General	16
9.2	Areas	16
9.3	Tests	17
9.3.1	General.....	17
9.3.2	Telephony	18
9.3.2.1	General	18
9.3.2.2	Scoring	18
9.3.3	Video streaming.....	19
9.3.3.1	General	19
9.3.3.2	Scoring	19
9.3.4	Data Testing.....	19
9.3.4.1	General	19
9.3.4.2	Scoring	20
9.3.5	Service Testing	20
9.3.5.1	General	20
9.3.5.2	Scoring	20
10	Statistical confidence and robustness	21
10.1	General	21
10.2	Influence of the derived scores on statistical confidence	21
10.3	Statistical confidence level estimation	22
10.3.1	General.....	22
10.3.2	Statistical analysis using a bootstrap resampling method	22
10.3.3	Interpretation of results	22
Annex A:	Example set of weighting factors, limits and thresholds	23
A.1	General	23
A.2	Area	23
A.2.1	Geographical divisions	23
A.2.1.1	General.....	23
A.2.1.2	City type	23
A.2.1.3	Road type.....	23
A.2.1.4	Complementary areas	23
A.3	Mobile services	24
A.4	Test metrics of mobile services	24
A.4.1	General	24
A.4.2	Telephony	24
A.4.3	Data Services.....	24
A.4.3.1	General.....	24
A.4.3.2	Video Streaming	24
A.4.3.3	Data Testing.....	25
A.4.3.4	Browsing.....	25
A.4.3.5	Social Media and Messaging	25
A.5	Example Calculation	26
Annex B:	Example set of weighting factors, limits and thresholds	27
B.1	General	27
B.2	Area	27
B.2.1	Geographical divisions	27
B.3	Mobile services	28
B.4	Test metrics of mobile services	28
B.4.1	Telephony.....	28
B.4.2	Data Services.....	28

B.4.2.1	General.....	28
B.4.2.2	Video Streaming	28
B.4.2.3	Data Testing	29
B.4.2.3.1	File Download (based on 3 MB File Size).....	29
B.4.2.3.2	File Upload (based on 1 MB File Size).....	29
B.4.2.3.3	File Download (based on 7 s Fixed Download Time).....	29
B.4.2.3.4	File Upload (based on 7 s Fixed Upload Time)	30
B.4.2.4	Browsing.....	30
B.4.2.4.1	Static Web Page	30
B.4.2.4.2	Dynamic Web Pages	30
B.5	Remarks on mapping functions.....	31
B.6	Example Calculation	32
History	33

Intellectual Property Rights

Essential patents

IPRs essential or potentially essential to normative deliverables may have been declared to ETSI. The information pertaining to these essential IPRs, if any, is publicly available for **ETSI members and non-members**, and can be found in ETSI SR 000 314: "*Intellectual Property Rights (IPRs); Essential, or potentially Essential, IPRs notified to ETSI in respect of ETSI standards*", which is available from the ETSI Secretariat. Latest updates are available on the ETSI Web server (<https://ipr.etsi.org/>).

Pursuant to the ETSI IPR Policy, no investigation, including IPR searches, has been carried out by ETSI. No guarantee can be given as to the existence of other IPRs not referenced in ETSI SR 000 314 (or the updates on the ETSI Web server) which are, or may be, or may become, essential to the present document.

Trademarks

The present document may include trademarks and/or tradenames which are asserted and/or registered by their owners. ETSI claims no ownership of these except for any which are indicated as being the property of ETSI, and conveys no right to use or reproduce any trademark and/or tradename. Mention of those trademarks in the present document does not constitute an endorsement by ETSI of products, services or organizations associated with those trademarks.

Foreword

This Technical Report (TR) has been produced by ETSI Technical Committee Speech and multimedia Transmission Quality (STQ).

Modal verbs terminology

In the present document "**should**", "**should not**", "**may**", "**need not**", "**will**", "**will not**", "**can**" and "**cannot**" are to be interpreted as described in clause 3.2 of the [ETSI Drafting Rules](#) (Verbal forms for the expression of provisions).

"**must**" and "**must not**" are **NOT** allowed in ETSI deliverables except when used in direct citation.

Introduction

Countrywide mobile network benchmarking and scoring campaigns published in the press enjoy great public interest and are of high importance for the operators of mobile networks. A first place score in press releases associated with such measurements is often used in the advertisements of the winning operator to boost their corporate identity. Though published results are often well documented, they are not always completely transparent about how the actual scoring has been achieved. Methods and underlying assumptions are mostly not described in detail.

The present document discusses the construction and methods of such a nationwide measurement campaign, with respect to the area and population to be covered, the collection and aggregation of the test results and the weighting of the various aspects tested. The applicability of the results of such a campaign, for inter country comparison purposes, is not covered in the present document.

Based on established methods and quality metrics, such as success ratio and setup times, the results of the data collected in the benchmarking are aggregated individually. The individual aggregated values are weighted and further aggregated for each application like telephony, video and data services. The application fields are then in turn weighted and aggregated over the different areas where the data is collected. Finally, calculation of an overall score or a joint score is performed.

The experienced quality of service varies over time so that the individual score of a particular throughput cannot be fixed once and for all. As well as the test metrics changing over time, so does the importance of the various services. The present document describes a typical set of tests that could be performed and a related evaluation criteria. In the annexes, actual real-world examples of weightings and score mapping parameters are given.

1 Scope

The present document describes the best practices for benchmarking of mobile networks. The goal of the benchmarking is to determine the best provider or operator for a designated area with respect of the services accessed with a mobile phone. The tests conducted are telephony, video streaming, data throughput and more interactive applications such as browsing, social media and messaging. This goal is achieved by executing benchmarking tests in designated test areas that represent or actually cover a major part of the users of mobile services. The results collected in the various areas are individually and collectively weighted and summarized into an overall score.

Due to the rapid development of the mobile technology and consumption habits of the users, the quality of experience of the users changes over time even when the objective to measure the quality of service does not change. The present document needs to keep up with those changes and does so by parameterizing the individual factors that contribute to the score.

2 References

2.1 Normative references

Normative references are not applicable in the present document.

2.2 Informative references

References are either specific (identified by date of publication and/or edition number or version number) or non-specific. For specific references, only the cited version applies. For non-specific references, the latest version of the referenced document (including any amendments) applies.

NOTE: While any hyperlinks included in this clause were valid at the time of publication, ETSI cannot guarantee their long term validity.

The following referenced documents are not necessary for the application of the present document but they assist the user with regard to a particular subject area.

- [i.1] ETSI TS 102 250-2: "Speech and multimedia Transmission Quality (STQ); QoS aspects for popular services in mobile networks; Part 2: Definition of Quality of Service parameters and their computation".
- [i.2] Void.
- [i.3] ETSI TR 102 505: "Speech and multimedia Transmission Quality (STQ); Development of a Reference Web page".
- [i.4] ETSI TR 101 578: "Speech and multimedia Transmission Quality (STQ); QoS aspects of TCP-based video services like YouTube™".
- [i.5] ETSI TR 102 678: "Speech and multimedia Transmission Quality (STQ); QoS Parameter Measurements based on fixed Data Transfer Times".
- [i.6] ETSI TR 103 138: "Speech and multimedia Transmission Quality (STQ); Speech samples and their use for QoS testing".
- [i.7] Recommendation ITU-T E.840: "Statistical framework for end-to-end network-performance benchmark scoring and ranking".
- [i.8] Recommendation ITU-T P.1401: "Methods, metrics and procedures for statistical evaluation, qualification and comparison of objective quality prediction models".
- [i.9] Recommendation ITU-T P.863: "Perceptual objective listening quality prediction".
- [i.10] Recommendation ITU-T P.863.1: "Application guide for Recommendation ITU-T P.863".

3 Definition of terms, symbols and abbreviations

3.1 Terms

For the purposes of the present document, the following terms apply:

live web page: web pages considered as dynamic content, content changes over time and some content might be different caused by the hosting server or the access network

static web page: web pages considered as static content, content stays constant over time and access network

3.2 Symbols

Void.

3.3 Abbreviations

For the purposes of the present document, the following abbreviations apply:

AMR	Adaptive Multi-Rate
API	Application Programming Interface
CDN	Content Delivery Network
CST	Call Setup Time
DL	DownLink
DNS	Domain Name System
EVS	Enhanced Voice Services
FB	FullBand
HD	High Definition
HTTP	HyperText Transfer Protocol
IP	Internet Protocol
ITU	International Telecommunication Union
ITU-T	International Telecommunication Union Telecommunication
KPI	Key Performance Indicator
MB	Megabyte
MOS	Mean Opinion Score
SMS	Short Messaging Service
TS	Technical Specification
UL	UpLink
VSSSR	Video Streaming Service Success Ratio
WB	WideBand

4 Governing Principles for Mobile Benchmarking

4.1 General

The accurate benchmarking and scoring of networks which cover large geographic areas requires careful consideration of a number of factors. These include the technology used, the extent of coverage offered, mobile device evolution, customer population distribution, network usage and tariff offerings. The following principles should be adhered to where possible to ensure that benchmarking scoring outcomes are always meaningful.

4.2 Fair Play

Benchmarking outcomes can be significantly influenced by specific targeting of test devices for superior performance. In such cases the results obtained no longer reflect the experience of a customer using that network. Steps should be taken to ensure that the measured results are truly representative of the real customer experience.

For example, if Operator A implements a special QoS construct specifically for the devices used to collect Benchmarking data, and Operator B does not, the results should not be compared for the purpose of drawing conclusions about the relative experience of customers on each network. The networks should not be compared for benchmarking purposes.

For example, if Vendor A implements a special functionality in their equipment/device software or firmware to recognize benchmark testing and boost performance, and Vendor B does not, the results may show one vendor to be superior to another for test cases no longer relevant to usual network usage. Vendor performance, from a customer perspective, can no longer be reliably compared.

4.3 Comparing networks with different coverage extents

Often networks are built with differing coverage objectives. Network rollout often varies between operators. This is often an important differentiator for customers making decisions about which network is best for them. Benchmarking should be performed in such a way that it highlights coverage differences in the results. From a scoring perspective, operators should never be penalized for providing coverage where other operators do not. In fact they should instead be rewarded in the scoring system. It should be the intention of any comprehensive mobile benchmark to include coverage comparison as a differentiating factor in the scoring.

For example, if Operator A offers significantly more geographic coverage than Operator B, Benchmarking data collection methodology and scoring should be such that this difference is always reflected in the scoring as a 'bonus' rather than a 'penalty' and the Benchmarking methodology should be such that this difference is measured. Failures occurring due to lack of coverage should always be included in scoring calculations and weighted appropriately to reflect the true customer experience.

4.4 Comparing networks with differing technology use

Network evolution and the adoption rate of new technologies often varies between operators. Benchmarking should be performed in such a way that it incorporates the use of the latest technology available. This is to reflect the network capability and customer experience available with the latest devices. Benchmark scoring should account for Operators who offer performance differentiation through early adoption of new technologies by way of a 'bonus' for such deployment.

For example, if operator A deploys 5G technology whilst operator B continues to deploy 4G technology, the benefits 5G technology offer to the customer experience should be captured in the Benchmarking data collection and scoring.

4.5 Test device selection

Mobile network benchmarking is performed mainly using drive testing. This relies heavily on the choice of test device(s). Care should be taken in the selection of such devices to ensure they do not favour one Operator's network over another in the results. The same devices may perform differently on two different networks depending on factors such as the antenna placement in the device for varying frequency bands, variations due to manufacturing tolerances, firmware version differences, modifications made to devices for metric data collection and device placement and mounting in the test vehicle.

4.6 Test server selection

Data tests are commonly performed to a test server or selected web page (or pages). The selection of such servers/sites can influence the benchmarking result. Test servers should be selected so they do not favour one network compared to another. Web pages should be selected such that they represent a cross section of pages commonly used by customers.

For example, if Operator A hosts the sever selected for 'ping' testing and the same server is also used to test Operator B, it is likely that performance levels for Operator B will be worse than those for Operator A due to the difference in latency to the selected server. This miss-represents the performance difference for this metric. Such situations should be avoided.

4.7 Test method transparency

Given the importance of the clear interpretation of benchmark results, all results should be accompanied by a declaration containing information about the following:

- 1) The scoring model/methodology used including all coefficients, targets and weightings.
- 2) The underlying KPI values as measured in the test.
- 3) The number of samples collected or number of tests performed for each KPI measured for each sub category.
- 4) The test methodology used including details of equipment setup, call sequences, test servers and web pages.
- 5) The areas/routes used for the data collection.
- 6) The device model and firmware version used for the data collection.
- 7) The tariff/data plan used for the data collection.

The intention of this is to provide the transparency required so that parties receiving the results are able to understand them fully. All factors required for this understanding should be provided.

4.8 Advice and best practice for web-page selection

Web page selection can impact on webpage load test results. To ensure a representative performance comparison can be made the following information and advice should be considered:

- For sufficient diversity and robustness of results, a minimum of 6 different pages is recommended to be considered for the scoring. It is good practice to measure more pages (e.g. 10), to retain enough diversity in case the dynamic behaviour requires to eliminate certain pages from the overall result.
- It is recommended to select pages according to their relevance to end customers. A popular ranking per country is given by Alexa Internet, Inc. (www.alexa.com). If possible, pages should be selected from Top 50 list, where an extension of that range is justifiable if not enough suitable pages exist within the Top 50.
- All pages should exceed a minimum size (e.g. 800 kB) to cover the minimum amount of data in case the download of a predefined data amount is used as success criteria. The page size needs to be observed on a daily basis throughout the measurements. In case of the severe size changes, a reaction may be needed.
- Internationally popular live pages and country dependent pages may be used in reasonable proportion (e.g. 10 live pages - 4 are common, 6 are country dependent).
- Ad blockers should not be used.
- A web-page selection that is hosted pre-dominantly by one CDN should be avoided.
- Websites of services that are predominantly accessed via a dedicated app on a smartphone should not be selected. For example, Facebook™, YouTube™ and similar websites/services are typically not accessed via a mobile browser and should therefore not be used as web-sites for HTTP Browsing tests in mobile benchmarking campaigns.
- No website should be selected that is a sub-page/site of another already selected website.
- No website should be selected where the content is legally suspicious or contains harming, racism or sexist content.

5 General Description

In the present document the benchmarking and scoring of networks over a large geographical area, e.g. entire countries in various modes and for diverse services provided by mobile networks is described. A comprehensive manner to compare the tested networks is to calculate an overall score per network based on the individual measurement results collected during a test campaign. The individual measurement results are aggregated using a weighted accumulation into an overall network score. This overall score finally allows the scoring of the tested networks. To arrive there, the weighted aggregation is done over several layers.

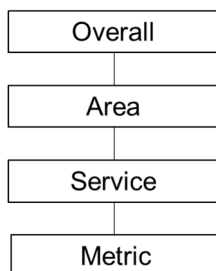


Figure 1: Aggregation layers

Weights are used for the aggregation of the different metrics, mobile services and areas to obtain the final score.

The accumulation of the measurements is done over several levels. The first or lowest layer consists of the measurement metrics for the services delivered over the mobile network. The services or applications considered are telephony, video, data transfer and services including browsing, social media and messaging. The metrics collected for one mobile service and a certain area are aggregated into an individual score for each metric; the scores of the metrics are then aggregated into an overall score of the mobile service.

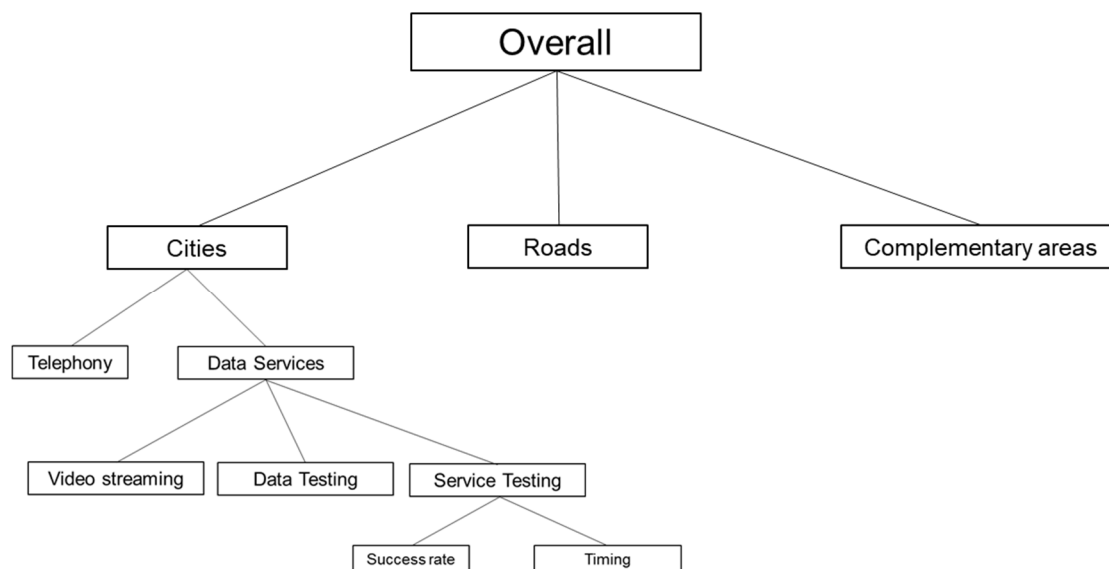


Figure 2: Aggregation over services and layers for a mobile network

In this aggregation, the metrics have a score weight according to the weight they were given for that particular mobile service. The scores for the individual mobile services are then in turn aggregated into a score for telephony and data services, and then together for the area they were collected in.

Finally the various areas are weighted and accumulated over the various areas covered in the measurement. The different areas can have further geographical subdivisions. The weighted aggregation of the areas results in an overall score that characterizes the network.

6 Test Areas

6.1 General

The choice of the areas to be tested are an important part of the test setup. In order to be representative the areas have to cover a majority of the population and main areas of mobile use; in case of limited countrywide coverage a representative proportion of the covered population. Drive testing is the method of choice but can be supplemented by walk testing in designated areas.

In the choice of areas and the distribution of time between individual subdivisions such as big cities and roads the geographical and topological properties of the respective country need to be considered. This may impair, to some extent, the comparability between countries. The aim should be that the chosen sites are appropriate for the respective country under test.

In order to be representative or to paint a more detailed picture, the areas of test such as cities and roads can be supplemented by measurements in trains and hot spot locations.

To maintain comparability, test areas that are not covered by all the networks under test need to be considered appropriately. In general, limiting the tests only to areas that are served by all networks is certainly the first choice, but in case important parts of the country and population would not be tested, the respective operator that does not cover these areas can be excluded from the countrywide testing or the limitations need to be included in the overall scoring.

The various areas need to be tested in an appropriate manner. Since some areas might not be accessible by drive testing, walk testing can be considered.

6.2 Geographical divisions

6.2.1 Cities

Cities are varying in size and density and the categorization of big, medium and small cities varies by country. The city size and importance is sometimes reflected in requirements set by the spectrum licensing authorities. The cities can be, but do not necessarily need to be, divided up into three categories, namely big cities, medium cities and small cities.

The big cities are defined as the major cities of a country from the population and commercial point of view. E.g. high rise buildings and high density of population are found in the big cities. Most of the hot spot areas are found in the big cities. Testing big cities means driving the main roads including tunnels and bridges.

Medium cities are smaller cities than the big cities with less inhabitants and less commercial importance. Occasionally they have high rise buildings and in general the density of the population is lower than in big cities.

Small cities or towns have fewer inhabitants than medium cities and have an even lower commercial importance.

The choice of the possible subdivision and distribution in defining city types is to reflect their relevance on the countrywide scale.

6.2.2 Roads

The highways are multi lane roads that can carry high traffic and connect big and medium cities of the test area. They are going across the country and have no intersections or traffic lights. Tests performed on city highways that are within a big or a medium city are counted in the results for cities rather than roads.

Main roads are roads that carry high traffic and connect cities of the test area. These roads may have traffic lights and intersections. The main roads that are driven within cities are counted for the cities.

Rural roads are roads that do not carry high traffic and connect medium and small cities. They can run through open landscape and can also cover dispersed settlements.

6.2.3 Complementary areas

Complementary tests, if appropriate, vary from country to country. E.g. trains and railways are an established locations for tests in countries with strong commuting or highly frequented intercity connections whilst in other countries trains can be disregarded.

Other hot spots of use such as train stations, airports, pedestrian zones, parks, stadiums or tourist attractions are locations frequented by users of mobile phones. Those areas are to be considered appropriately.

7 User Profiles

Different users have different requirements and expectations with regards to mobile services. These expectations are the basis of what is perceived as excellent, good or poor. In addition, the type of service that is requested might differ between different user groups whom each put a different emphasis on the various service aspects like telephony, video, data or other social media. These groups can be assigned different subscription profiles, however, for the purpose of a network comparison, the best or highest commercially available profile yields the results that represent the performance of the network in best fashion. Using standard or budget profiles may produce interesting insights in the services received by the respective subscriber but are not in the position to assess what the network is capable of.

8 Test Metrics

8.1 Introduction

The test metrics are, with a few exceptions, generally defined in ETSI TS 102 250-2 [i.1] therefore the tests are whenever possible referenced to that document. In ETSI TS 102 250-2 [i.1] a success ratio is in most cases a two-step metric divided into successful access and successful conclusion.

These can be by weighted calculation aggregated into a single value, possibly even incorporating additional metrics or criteria that are decisive for the user's perception of a working service. This applies to all occurrences of Success Ratio in clauses 8.2 to 8.5.

The video, data throughput and service testing are often summarized as data tests as opposed to telephony tests, this separation is not excluded.

8.2 Telephony

8.2.1 General

Telephony tests are tests with a fixed call length where two terminals, either both mobile or one landline and one mobile call each other. Landline connections usually do not support new higher codecs such as AMR WB or EVS. In order to measure these codecs mobile to mobile circuit switched calls are necessary and at times even VoLTE calls over packet switched are needed. To consider unsustainable quality in a call, for a low speech quality score (e.g. MOS < 1,6) or silent periods for consecutive measurement samples (e.g. > 20 s), the call can be counted as unsustainable, and as an unsuccessful call or treated by a separate indicator.

8.2.2 Telephony Success Ratio

The success rate of the voice service independent of access or relay technology is the Telephony non Accessibility and the Telephony Cut-Off Call Ratio in ETSI TS 102 250-2 [i.1], clauses 6.6.1 and 6.6.5. For the purpose of the present document, the voice over LTE (VoLTE) service is treated as a telephony service.

8.2.3 Setup Time

The setup time for voice calls is defined in ETSI TS 102 250-2 [i.1], clause 6.6.2. It starts with the initiation of the call and ends when the alerting of the called side is indicated. Alternatively, the time when the acceptance or successful setup of the call is signalled to the user can be used as the end trigger.

8.2.4 Listening Quality

The value is calculated on a per sample basis as described in ETSI TS 102 250-2 [i.1], clause 6.6.4 where Recommendation ITU-T P.863 [i.9] in FB mode needs to be used. The measurement is set up according to ETSI TR 103 138 [i.6] and Recommendation ITU-T P.863.1 [i.10].

8.3 Video Testing

8.3.1 General

Video testing is in the standard case IP based video streaming. Video streaming quality of service aspects can be found in ETSI TS 102 250-2 [i.1] and in ETSI TR 101 578 [i.4]. For the purposes of the present document the Smartphone app based testing as in Figure 1 of ETSI TR 101 578 [i.4] is used. In order to collect details of the transport and reproduction, the length of the observation period of the video should reflect the relevant delivery mechanisms and the typical usage profile of a mobile user.

8.3.2 Video Streaming Service Success Ratio

The video streaming success ratio is the end-to-end success ratio of the requested video stream. It starts with the request of the video and ends with the end of the playout. This is derived from the metrics in ETSI TR 101 578 [i.4] as a combination of Video Access Failure Ratio and Video Playout Cut-off Ratio.

8.3.3 Setup Time

The setup time is the time from stream request to the display of the first picture and start of playout. This is Video Access Time from ETSI TR 101 578 [i.4].

8.3.4 Video Quality

The quality of a video reproduction is determined by freezing, frame-rate, resolution and compression depth and scheme by the codec. Freezing is most common and annoying impairment experienced by the user. The handling of freezes is described in clause 4.5.4 in ETSI TR 101 578 [i.4].

A comprehensive measure for the perceived quality that combines the impact of the above mentioned parameters is the mean opinion score (MOS) scale and is done according to clause 6.5.8 in ETSI TS 102 250-2 [i.1]. In the case of video streaming with a respective app on the smartphone an encrypted stream and a range of different resolutions (up to HD) is expected. The Video Quality parameter in ETSI TR 101 578 [i.4] reflects such measure. In addition to this, Video Freezing Time proportion in ETSI TR 101 578 [i.4] provides an insight about the proportion of the accumulated video freezing duration in relation to the actual video playout duration.

8.4 Data Testing

8.4.1 General

For data testing the throughput bandwidth for the user is tested. This is done by downloading and uploading incompressible files over HTTP. In clause 6.8 of ETSI TS 102 250-2 [i.1] the up and download of entire files is described. The description of an upload and download using fixed duration is in clause 5.2 of ETSI TR 102 678 [i.5]. Both approaches can be used, either alone or combined, for the purpose of evaluating throughput bandwidth.

8.4.2 Success Ratio

The determination of the success ratio for HTTP uploads and downloads is included in clause 6.8 of ETSI TS 102 250-2 [i.1] and in clause 5.2 of ETSI TR 102 678 [i.5].

8.4.3 Throughput

The determination of the mean data rate or throughput for HTTP uploads and downloads is included in clause 6.8 of ETSI TS 102 250-2 [i.1] and in clause 5.2 of ETSI TR 102 678 [i.5]. The best view of the actual download bandwidth is provided by a multi-threaded HTTP download test.

8.5 Services Testing

8.5.1 General

Besides the browsing of web pages as in clause 6.8 of ETSI TS 102 250-2 [i.1], services like social media and messaging systems (SMS is not considered) are not described or standardized for mobile testing. Some overall interesting aspects of all of these services are the success ratio and the duration or timing of the interaction.

8.5.2 Services

8.5.2.1 Browsing

For web browsing tests web pages are accessed and downloaded. These pages can be static like the Kepler page ETSI TR 102 505 [i.3] and - preferred - popular dynamic pages. The respective browsing metrics are in clause 6.8 of ETSI TS 102 250-2 [i.1]. For dynamic web pages a success criterion can be defined by the time it takes until a predefined data volume of the overall page session is received.

8.5.2.2 Social Media

For social media like Facebook™ and Instagram™, an action such as posting of pictures, text and video is the typical activity that is tested.

These are the usual activities where the user interacts with the media application. Popular social media vary in their popularity over time and across countries, therefore the list of services and their weight in the calculation can change. Since at the time of publication of the present document there are no standardized metrics for the use of these services over mobile networks, the metrics cannot be referenced to any document. In cases where interfaces (API) exist to those applications, these can be used to test the respective service.

It should be highlighted that the perceived user experience depends on the tested service platform in addition to the mobile network performance, because the observed timings inevitably include processing time on the service platform. As such, it depends on the focus of the testing activity as to whether inclusion of such services is useful or not.

8.5.2.3 Messaging

Sending a text message, line and measuring the delivery time and success rate is a convenient way to characterize the perceived quality of service. For these services, delivered over a mobile network, no standard for quality of service metrics exists.

Although there are no standardized metrics for social media and messaging services delivered over mobile network, metrics based on legacy messaging services can be established.

8.5.3 Success Ratio

In general the successful conclusion of an activity is to be measured in social media and messaging. The number of successful trials versus the number of trials is the success ratio.

$$\text{Service Success Ratio [\%]} = \frac{\text{number of successful activities}}{\text{number of trials}} \times 100$$

An activity starts with triggering an action on the device by e.g. pushing a button to send a text message, to open a Facebook™ profile, posting a picture on Instagram™ or opening a web page. The activity is successful when the application indicates a confirmation that the triggered process is successfully concluded. This can be done, for example, by a graphic indicator like a check or by other means.

8.5.4 Timings

The duration of a social media or messaging activity is the time between triggering the activity and the indication of the successful conclusion of it. In the case of browsing, social media and for messaging, it is the time until confirmation of successful reception is indicated.

$$\text{Service activity Duration [s]} = t_{\text{end}} - t_{\text{start}}$$

The timing, the duration from the initiation to the successful conclusion of a test depends to a significant extent on the performance of the underlying web service. However, these factors are the same for all networks under test.

9 Weighting

9.1 General

In order to achieve an overall score the individual test results for the various areas have to be given a weight. This weight is the importance with which the result enters into the overall valuation of the testing. The weighting of the results is done on each level of aggregation (Figure 1).

In the following clauses the general method of weighting and the individual measures are described.

Example values that are used in practice for the weighting of the areas and tests as well as actual values for the upper and lower limits of the target ranges are presented in the annexes to the present document.

9.2 Areas

For an area, all regional, daytime, geographical or morphologic categories are considered, where the scoring method is applied before further aggregating to an overall score. These different categories that are measured have a combined weight of 100 %, in case e.g. there are no complementary areas, cities and roads have alone a combined weight of 100 %. In case the areas have further subdivision these areas are individually weighted and then make up 100 % of the next level e.g. if city category is subdivided into big cities and small cities. These two subareas add up to 100 % representing the whole weight of cities.

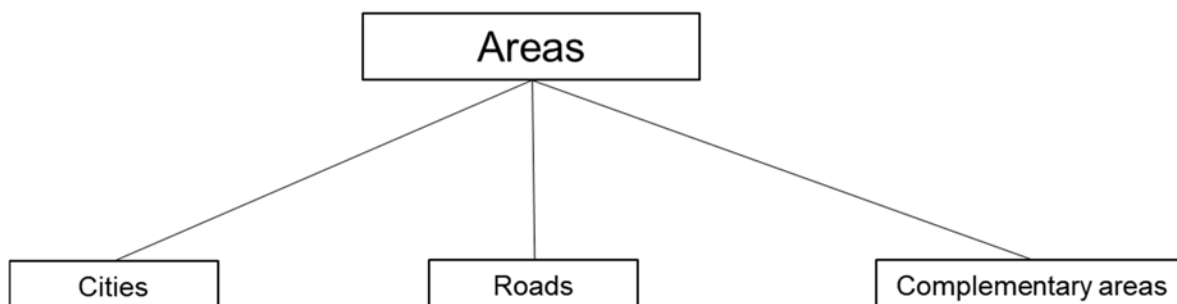


Figure 3: Examples of areas

The timing of tests can have an impact on the perceived performance on the network. Despite this, the weighing of different times is not considered part of best practice. It is advisable, however, that measurements are reasonably spread over the different times of day, e.g. to not deliberately exclude busy hours.

9.3 Tests

9.3.1 General

Each test is multi layered in nature. The upper layer provides the overall score of the mobile service tests, which is calculated from the weighted scores of the test scenarios for telephony and data services. The two scenarios have the combined weight of 100 %. The data services in turn consist of video streaming, data testing and service testing. The three types have also the combined weight of 100 %. The weight of the individual test types can be determined according to the intended user profile.

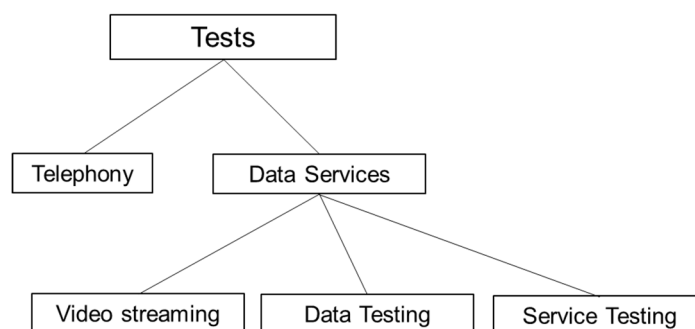


Figure 4: Service types for testing

The test metrics are evaluated as aggregated values. While the success rate is aggregated already, for most of the other values such as listening quality, throughput, setup time, and duration etc. the average is taken into account. These individual metrics have a minimum and a maximum value. However, the metrics do also have a bad limit to saturation area in which the experience of the customer does not deteriorate significantly and a good limit to saturation area above which the customers experience does not improve further. The average values are expected to be between the good and the bad limits.

The scoring of the individual aggregates can be increasing or decreasing. If the score rises with the value then it is an increasing value score. Starting from the minimum below the bad limit the value score is at 0 %. Between the bad and the good limit, the score is increasing to 100 %. In the saturation area between the good limit and the maximum, the value score stays at 100 %.

If the value score rises with the decreasing value, then it is a decreasing value score. In the saturation area between the good limit and the minimum value, the score is stable at 100 %. Between the good and the bad limit, the score decreases to 0 % and stays there above the bad limit.



Figure 5: Weighting function

The general formula is:

$$\text{Score} = \frac{\text{value} - \text{Bad limit}}{\text{Good limit} - \text{Bad limit}} \times \text{weight}$$

NOTE: In the case of a decreasing value score two things are expected:

- 1) the bad limit will be a higher numerical value than the good limit;
- 2) the resulting negative/negative calculation is expected and produces a positive result.

The scores for the average values are calculated between two limits. The negative impact of poor performance or the positive impact of excellent performance can be underrepresented by taking only scored averages into account. The aspects of the distribution of the results need to be taken into account. It is therefore useful to introduce limits for poor and excellent performance and calculate the percentage or percentile of results within these limits. In order to boost superior performance, an extra bonus can be applied, similarly the achievement of minimum performances can be awarded too.

In the given graph in Figure 5, a linear function is given to illustrate the general method, although not all service measurements are perceived in a linear manner. A non-linear relationship may occur where an increase or decrease in metric value at one end of the scoring interval is perceived to have a much larger or smaller effect than a similar increase or decrease at the other end and hence should have a higher or lower percentage of the available score. In this case non-linear functions may be applied to determine the score value. There are a number of functions which could be utilized to measure the non-linear score with square root or logistic sigmoid function being typically and widely used at this time. A working example for the application of a non-linear weighting is given in Annex B of the present document.

9.3.2 Telephony

9.3.2.1 General

The telephony service has three major aspects: overall success rate, setup time and listening quality (MOS). These three values enter into the calculation of the overall score of the telephony service. The individual aspects can then in turn be weighted individually for the calculation of the overall telephony score. These factors have a combined weight of 100 %.

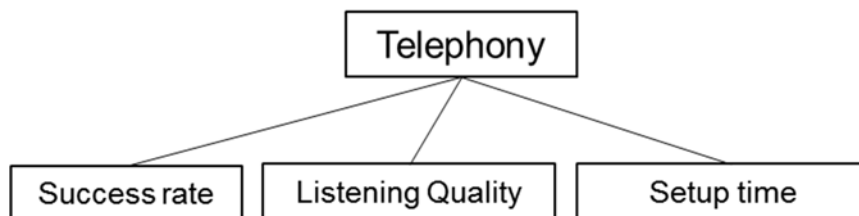


Figure 6: Contributing dimensions to telephony

9.3.2.2 Scoring

The higher the call success rate and the MOS value the better is the experience of the user; the two values have an increasing value score. While the longer the call setup time is, the worse the experience of the user is; the call setup time has a decreasing score value. In addition to the consideration of the average values, extra bonus for excellent listening quality or very short call setup-times can be given, the same as extra bonus for the reduction of very bad experiences as very low MOS scores or very long call setup times. As examples, the 10th percentile of CST can be used for awarding very short CST and the 90th percentile for awarding excellent listening quality. To award the absence of negative experiences the 90th percentile of CST can be considered and a ratio of MOS < threshold (e.g. MOS < 1,6 or MOS < 2,2). A lower parameter value would lead to a higher score in these cases.

For telephony the following factors with example thresholds and percentiles can be taken into account:

- Call Setup Success Ratio.
- Call Drop Ratio.
- MOS.
- MOS < low MOS threshold.

- 90th percentile of MOS.
- Call Setup Time.
- Call Setup Time > long setup time threshold.
- 90th percentile of Call Setup Time.

9.3.3 Video streaming

9.3.3.1 General

The main aspects of the video streaming are the video streaming service success ratio, setup time (video access time) and the visual quality. These three factors are combined to give the video streaming score together with extra bonus for superior performance values for selected metrics. All factors have the combined weight of 100 %.

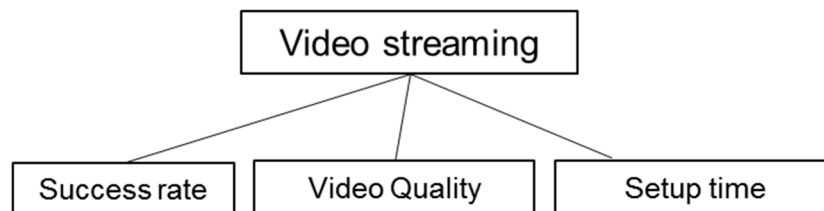


Figure 7: Contributing dimensions to video streaming

9.3.3.2 Scoring

The higher the video streaming service success ratio and the video quality MOS value the better the experience of the user is. The two values have an increasing value score, while the longer the time to first picture display and video start picture (that is video access time) is the worse is the experience of the user is; the setup time has a decreasing score value. The negative impact of bad video quality is represented by e.g. 10 percentile of video quality MOS between two limits and the negative impact of long setup times is represented by the percentage of video access time is above e.g. 10 s. Any impact of particularly good performance is not taken into account for video streaming since there is a very high proportion of HD expected that does not leave much headroom for technical improvement that can be rewarded with a bonus. For video streaming, the following factors with the proposed thresholds and percentiles can be taken into account:

- Video Success Ratio.
- Video Quality.
- Video Freezing Time proportion.
- Resolution.
- 10th percentile of Video Quality.
- Video Access Time > long setup time threshold.

The streaming success ratio together with the quality measures such as MOS, freezing, resolution and video access time can be combined to define a composite success criterion, with minimum requirements for the quality metrics. In this case, only video sessions are scored for quality aspects, which succeed in the composite success criterion.

9.3.4 Data Testing

9.3.4.1 General

The main aspects of data testing are the success rate and the data rate or throughput. These two factors are combined to produce the data testing score. These have a combined weight of 100 %.

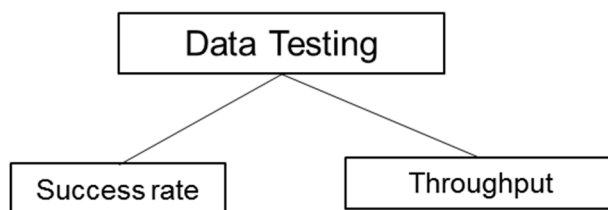


Figure 8: Contributing dimensions to data transfer testing

9.3.4.2 Scoring

The higher the success rate and the throughput value the better the experience of the user is; the two values have an increasing value score. The negative impact of low throughput values is represented by the e.g. 10th percentile and the positive impact of high throughput is represented by e.g. 90th percentile. The Average Session Duration has a decreasing value score.

For data testing, the following factors with the proposed thresholds and percentiles can be taken into account. These apply for both uplink and downlink:

- Transfer Success Ratio.
- Average Session Duration.
- Average throughput.
- 10th percentile of (low) throughput.
- 90th percentile of (high) throughput.

9.3.5 Service Testing

9.3.5.1 General

The main aspects of the services are the success rate and the timing or duration. These two factors are combined to determine the service score for browsing, social media and messaging. These aspects have a combined weight of 100 %.

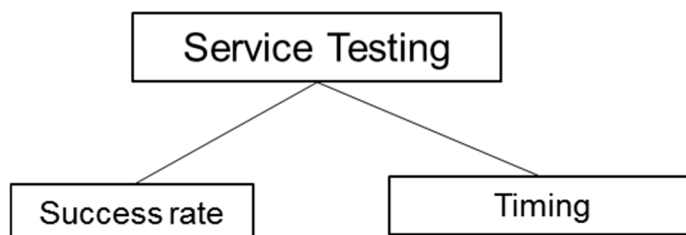


Figure 9: Contributing dimensions to data service testing

9.3.5.2 Scoring

The higher the success rate the better the experience of the user. The success rate value has an increasing score value. The timing or duration of an activity such as browsing or posting follows a decreasing function. The longer it takes the worse is the experience is. The negative impact of long activity duration is represented by the percentage of times above a long duration threshold.

For service testing, the following factors are taken into account:

- Activity Success Ratio.
- Average Duration.
- Activity Duration > long duration threshold.

For social media and messaging services only the activity duration samples of sending text message and sharing single picture can be taken into the factor calculation.

For Social Media and Messaging service testing, the following factors can be taken into account:

- Activity Success Ratio.
- Average Duration.
- Activity Duration > long duration threshold.

10 Statistical confidence and robustness

10.1 General

When comparing mobile networks, through benchmark metrics and scoring means, it is important that the statistical significance of the outcome is considered. This is especially true when using the results to conclude if one network is more preferable than another. The following clauses outline important considerations and provide one method to evaluate the validity of conclusions from a statistical viewpoint.

10.2 Influence of the derived scores on statistical confidence

The performance of a given network is estimated based on a set of measurements. These measurements are collected in certain geographical locations, at certain times on the dates selected to represent the general experience of customers under real-field load conditions. The measurements from any measurement campaign only represent a subset of measurements from the overall population.

The measured attributes and the derived metrics/scoring are subject to uncertainty. Two questions arise. How closely do the results derived from the measurements represent the performance experienced by the entire population? How repeatable are the results, it means how sensitive are the results to changes in the measurement points within the same basic population?

In general, the larger the sample set the less the uncertainty of the results, and the better representation of the population distribution of measurement results. Therefore, it is necessary that the sample set is selected to include samples from the various different environments that exist for the overall population. The uncertainty of the indicators and the derived score requires statistical analysis and is usually described by confidence intervals.

It should be mentioned that a particular measured attribute will have an individual confidence interval which may result in a stable average with more or less samples than another attribute. Contributors with low confidence typically drive and decrease the confidence of a final aggregated score. This should be considered when defining and scaling measurement campaigns. Low confidence stands for larger statistical confidence levels.

Success or failure ratios, in particular, require a sufficient amount of measurements to support a confident conclusion. For example, in a collection of 100 calls, one dropped call more or less will lead to a change in the Call Drop Ratio of 1 %, in case of 1 000 calls the Call Drop Ratio will change by 0,1 % by a single dropped call. The resolution of the Call Drop Ratio is defined by the number of measurements from which it is derived. If the actual Call Drop Ratio is within the range limited by the measurement resolution, deriving a reliable representation is difficult. In such cases measurement sample volumes should be increased. The minimum number and the targeted confidence interval for benchmarking campaigns highly depends on the purpose of the campaign.

10.3 Statistical confidence level estimation

10.3.1 General

This clause presents a pragmatic, empirical approach to assessing basic statistical properties of network benchmarking score results and deriving confidence levels. It will answer these main questions:

- How close can one expect the unknown results of the basic population of the observation area to be to the results and score of the obtained measurement result?
- How repeatable are the obtained measurement results and score when using different measurement points of the same basic population?

The following method describes an established method for deriving a confidence interval and other statistical metrics of the scoring result. The derived confidence interval is based on statistical evaluation of the measurement results and is a pure statistical representation of the measurement samples. It does not reflect the significance in human perception.

10.3.2 Statistical analysis using a bootstrap resampling method

The presented method describes a probability density (and distribution function) and the related statistical properties (average, standard deviation, confidence intervals) of the achieved score of each competing network. This helps to interpret results and establish a means of determining the confidence of an achieved score from the collected measurement data.

The probability density function of a scoring result for a network or an area can be derived following a bootstrapping approach according to this pseudo-algorithm:

- 1) For each measurement contributing to the score, use re-sampling from empirical distribution (i.e. set of measured values) to generate a 'bootstrap sample' of equal sample size as the set of measured values.
- 2) Calculate relevant statistics (i.e. aggregated QoS Parameters, score) from the bootstrap sample and map the result to the score domain.
- 3) Repeat steps 1) and 2) for a sufficiently large number of times (N) to estimate sufficiently the distribution of the QoS parameters or the aggregated score.
- 4) Assess the statistical properties of the result, e.g. confidence intervals, standard deviation, etc.

The respective algorithm is performed for each aggregated QoS Parameter. It is advisable that the re-sampling, though random in nature, be executed in a way that preserves (within the scope of one bootstrap subsample) the correlation between QoS parameters originating from the same service session.

Assuming bootstraps of size $M = 1\,000$ are being created for each aggregated QoS Parameter and there are $N = 10\,000$ of those bootstrap samples, the resulting distribution of scores can be interpreted as the results of 10 000 hypothetical testing campaigns carried out under similar conditions. Having obtained these hypothetical results (i.e. distribution of scores) for a network or an area the confidence interval of the score or even an individual contributor can be estimated. Further information can be found in Recommendation ITU-T E.840 [i.7] and Recommendation ITU-T P.1401 [i.8].

10.3.3 Interpretation of results

The method described in the previous clause allows assessing empirically the statistical variation of the score and its statistical confidence interval. The statistical confidence interval also allows analysis of whether the difference between two scores (e.g. two areas, two networks, two different time ranges) is significant by certain probability (e.g. 95 %). The significance of such a difference is purely based on the statistical evaluation and will not give any indication of the significance in human perception of the networks performance.

Furthermore, applied significance levels and comparisons depend on the purpose of the measurement and are not generalized or recommended in the present document.

Annex A: Example set of weighting factors, limits and thresholds

A.1 General

This annex provides a first example which represents a best practise at the time of release of the present document. The information here is intended to provide an illustration of how to practically apply network benchmarking and scoring as described in the body of the present document. In this regard it identifies example weights, limits and thresholds that could be applied to areas and mobile services as well as providing example worked network scoring calculations.

A.2 Area

A.2.1 Geographical divisions

A.2.1.1 General

The three areas can be weighted in the following manner.

Area	Weight
Cities	50 %
Roads	40 %
Complementary areas	10 %

A.2.1.2 City type

In case of three subdivisions of the cities a possible weighting is as follows.

City Type	Weight
Big cities	60 %
Medium cities	30 %
Small cities	10 %

A.2.1.3 Road type

The three types of roads can be weighted in the following manner.

Road Type	Weight
Highways	60 %
Main Roads	30 %
Rural Roads	10 %

A.2.1.4 Complementary areas

The two general walk tests can be weighted in the following manner.

Type	Weight
Trains	40 %
Hotspots (train stations, airports, pedestrian zones, parks, stadiums or tourist attractions)	60 %

A.3 Mobile services

Service Type	Weight
Telephony	40 %
Data Services	60 %

A.4 Test metrics of mobile services

A.4.1 General

For all the considered service types, the basic idea is to rate three aspects of the service, where possible:

- Service availability/retainability (e.g. success ratio, drop ratio).
- Service access time (e.g. video access time or call setup time).
- Quality of the media transfer (e.g. listening quality, time to present the webpage content).

Each of the aspects is described and scored by one or more quality indicators. Usually, there is one indicator for rating the average performance (e.g. average call setup time), other indicators rate superior or low performance (e.g. 10th and 90th percentile or the ratio of tests exceeding a certain threshold).

A.4.2 Telephony

Factor	Bad limit	Good limit	Weight
Call Setup Success Ratio	85,00 %	100,00 %	31,25 %
Call Drop Ratio	10,00 %	0,00 %	37,50 %
MOS	2,00	4,30	4,38 %
MOS < 1,6	10,00 %	0,00 %	5,62 %
90 th percentile of MOS	4,00	4,75	2,50 %
Call Setup Time [s]	12,00	4,50	6,25 %
Call Setup Time > 15 s	3,00 %	0,00 %	8,75 %
90 th percentile of Call Setup Time	8,00	4,00	3,75 %

A.4.3 Data Services

A.4.3.1 General

Data Service Type	Weight
Video Streaming	22 %
Data Testing	25 %
Browsing	38 %
Social Media and Messaging	15 %

A.4.3.2 Video Streaming

Factor	Bad limit	Good limit	Weight
Video Streaming Service Success Ratio (VSSSR)	80,0 %	100,0 %	58,00 %
Video Quality MOS	3,0	4,5	16,50 %
10 percentile of MOS	2,0	4,0	16,50 %
Video access time [s]	7,0	2,0	4,50 %
Video access time > 10 s	5,0 %	0,0 %	4,50 %

A.4.3.3 Data Testing

Factor	Bad limit	Good limit	Weight
Transfer Success Ratio DL (e.g. 5 MB)	80 %	100 %	11,00 %
Average throughput DL [Mbit/s]	1,0	100,0	14,00 %
10 th percentile of (low) throughput DL [Mbit/s]	1,0	40,0	18,00 %
90 th percentile of (high) throughput DL Mbit/s]	10,0	240	7,00 %
Transfer Success Ratio UL (e.g. 2 MB)	80 %	100 %	11,00 %
Average throughput UL [Mbit/s]	0,5	50,0	14,00 %
10 th percentile of (low) throughput UL [Mbit/s]	0,5	30,0	18,00 %
90 th percentile of (high) throughput UL [Mbit/s]	5,0	100,0	7,00 %

A.4.3.4 Browsing

Factor	Bad limit	Good limit	Weight
Activity Success Ratio	80,0 %	100,0 %	66,67 %
Average Duration [s]	6,0	1,0	28,57 %
Activity Duration > 6 s	15,00 %	0,00 %	4,76 %

A.4.3.5 Social Media and Messaging

Factor	Bad limit	Good limit	Weight
Activity Success Ratio	80,0 %	100,0 %	66,67 %
Average Duration [s]	15,0	3,0	28,57 %
Activity Duration > 15 s	5,00 %	0,00 %	4,76 %

A.5 Example Calculation

	Bad limit	Good limit	Weight in service	Weight in Data or Telephony	Weight of Data or Telephony	Example results City	RAW =MIN(MAX((F-A)/(B-A)) *100 ; 0);100)	Score = RAW * C * D * E	Example results Rural	RAW =MIN(MAX((G-A)/(B-A)) *100 ; 0);100)	Score = RAW * C * D * E
Column	A	B	C	D	E	F			G		
Telephony											
Call Setup Success Ratio	85%	100%	31.25%	100%	40%	96%	73.3	9.167	90%	33.3	4.167
Call Drop Ratio	10%	0%	37.50%	100%	40%	2.1%	79.0	11.850	3.2%	68.0	10.200
MOS	2	4.3	4.38%	100%	40%	3.4	60.9	1.066	3.3	56.5	0.990
MOS < 1,6	10%	0%	5.62%	100%	40%	2.85%	71.5	1.607	4.60%	54.0	1.214
90th percentile of MOS	4	4.75	2.50%	100%	40%	3.86	0.0	0.000	3.74	0.0	0.000
Call Setup Time [s]	12	4.50	6.25%	100%	40%	4.4	100.0	2.500	6.4	74.7	1.867
Call Setup Time > 15 s	3%	0%	8.75%	100%	40%	1.1%	63.3	2.217	2.4%	20.0	0.700
90th percentile of Call Setup Time	8	4	3.75%	100%	40%	3.6	100.0	1.500	3.7	100.0	1.500
Video Streaming											
Streaming Success Ratio	80%	100%	58.00%	22%	60%	89%	45.0	3.445	81%	5.0	0.383
Video Quality MOS	3	4.5	16.50%	22%	60%	3.7	46.7	1.016	3.3	20.0	0.436
10 percentile of MOS	2	4	16.50%	22%	60%	2.7	35.0	0.762	2.2	10.0	0.218
Video Access Time [s]	7	2	4.50%	22%	60%	3.2	76.0	0.451	3.9	62.0	0.368
Video Access Time > 10 s	5%	0%	4.50%	22%	60%	2.8%	44.0	0.261	4.3%	14.0	0.083
Data Testing											
Transfer Success Ratio DL (e.g. 5MB)	80%	100%	11.00%	25%	60%	96%	80.0	1.320	86%	30.0	0.495
Average throughput DL [Mbit/s]	1	100	14.00%	25%	60%	44.4	43.8	0.921	25.7	24.9	0.524
10th percentile of (low) throughput DL [Mbit/s]	1	40	18.00%	25%	60%	12.8	30.3	0.817	9.5	21.8	0.588
90th percentile of (high) throughput DL [Mbit/s]	10	240	7.00%	25%	60%	78.4	29.7	0.312	52.8	18.6	0.195
Transfer Success Ratio UL (e.g. 2MB)	80%	100%	11.00%	25%	60%	92%	60.0	0.990	95%	75.0	1.238
Average throughput UL [Mbit/s]	0.5	50	14.00%	25%	60%	9.1	17.4	0.365	4.2	7.5	0.157
10th percentile of (low) throughput UL [Mbit/s]	0.5	30	18.00%	25%	60%	0.98	1.6	0.044	0.76	0.9	0.024
90th percentile of (high) throughput UL [Mbit/s]	5	100	7.00%	25%	60%	14.7	10.2	0.107	10.6	5.9	0.062
Browsing											
Activity Success Ratio	80%	100%	66.67%	38%	60%	96%	80.0	12.161	95%	75.0	11.401
Average Duration [s]	6	1	28.57%	38%	60%	3.6	48.0	3.127	3.9	42.0	2.736
Activity Duration > 6 s	15%	0%	4.76%	38%	60%	8%	44.7	0.485	11%	26.7	0.289
Social Media and Messaging											
Activity Success Ratio	80%	100%	66.67%	15%	60%	87%	35.0	2.100	86%	30.0	1.800
Average Duration [s]	15	3	28.57%	15%	60%	10.4	38.3	0.986	12.1	24.2	0.621
Activity Duration > 15 s	5%	0%	4.76%	15%	60%	3%	36.0	0.154	13%	0.0	0.000
							Sum (City)	59.73		Sum (Rural)	42.26
							Weight	60%		Weight	40%
							Overall Score (60% City + 40% Rural)				52.74

Annex B: Example set of weighting factors, limits and thresholds

B.1 General

This annex provides a second example which represents a best practise at the time of release of the present document. The information here is intended to provide an illustration of how to practically apply network benchmarking and scoring as described in the body of the present document. In this regard it identifies example weights, limits and thresholds that could be applied to areas and mobile services as well as providing example worked network scoring calculations.

B.2 Area

B.2.1 Geographical divisions

The areas could be weighted and subdivided in the following manner.

EXAMPLE 1:

Area Type	Cumulative Area Type Weight	Area	Subdivision	Weight
Cities	80 %	Cities	Big and Medium	45 %
		Cities	Towns	20 %
		Complementary Areas	Hotspots	15 %
Outside Cities	20 %	Roads	n/a	12,5 %
		Complementary Areas	Railways	7,5 %

The Area Type is not a dimension introduced in the main part of the present document and serves merely to illustrate how the introduced Areas and their subdivisions are logically grouped.

This allows for alternatives which keep the high-level distribution between Area Types if, depending on the scope of the exercise, other combinations are possible where complementary areas are not in scope.

Two further examples are shown below.

EXAMPLE 2:

Area Type	Cumulative Area Type Weight	Area	Subdivision	Weight
Cities	80 %	Cities	Big and Medium	45 %
		Cities	Towns	20 %
		Complementary Areas	Hotspots	15 %
Outside Cities	20 %	Roads	n/a	20 %

EXAMPLE 3:

Area Type	Cumulative Area Type Weight	Area	Subdivision	Weight
Cities	80 %	Cities	Big and Medium	60 %
		Cities	Towns	20 %
Outside Cities	20 %	Roads	n/a	20 %

B.3 Mobile services

Service Type	Weight
Telephony	40 %
Data Services	60 %

B.4 Test metrics of mobile services

B.4.1 Telephony

Factor	Lower limit		Upper limit	Weight
	Cities	Outside Cities		
Composite Call Success Criterion combining: <ul style="list-style-type: none"> – Call Setup Success Ratio – Call Setup Time < 15 s – Inverse of Call Drop Ratio – Ratio of calls with no 2 consecutive speech samples < 1,3 MOS 	87 %	80 %	100 %	70 %
Average of MOS across all Samples	2,5 MOS		4,4 MOS	15 %
Average of Call Setup Time [s]	2,5 s		8 s	3 %
90 th percentile of Call Setup Time [s]	3 s		10 s	12 %

B.4.2 Data Services

B.4.2.1 General

Data Service Type	Weight	Subdivision	Subdivision Weight
Video Streaming	22 %		
Data Testing	56 %	Fixed Size File DL	14 %
		Fixed Size File UL	14 %
		Fixed Duration File DL	14 %
		Fixed Duration File UL	14 %
Browsing	22 %	Static Page (Kepler)	16 %
		Dynamic Pages	6 %

B.4.2.2 Video Streaming

Factor	Lower limit		Upper limit		Weight
	Cities	Outside cities	Cities	Outside cities	
Composite Session Success Criterion combining: <ul style="list-style-type: none"> • Success to access and stream the video (Video Streaming Service Success Ratio, VSSSR) • Absence of freezing beyond tolerable level 	87,0 %	80,0 %	100,0 %		55,0 %
Average Video Resolution [p]	480 p		1 080 p		25,0 %
Ratio of videos with no freezing [%]	95,0 %	90,0 %	100,0 %		15,0 %
Video access time [s]	1,5 s		4,5 s	5,0 s	5,0 %

B.4.2.3 Data Testing

B.4.2.3.1 File Download (based on 3 MB File Size)

Factor	Lower limit		Upper limit		Weight
	Cities	Outside cities	Cities	Outside cities	
Composite Session Success Criterion combining: <ul style="list-style-type: none"> Success to access and download the file Achievement of min. throughput of 384 kbit/s 	87,0 %	80,0 %	100,0 %		60,0 %
Average Download Session Duration [s]	1,0 s	1,5 s	10,0 s		12,0 %
10 th percentile of Download throughput [kbit/s]	384 kbit/s		27 500 kbit/s	11 000 kbit/s	20,0 %
90 th percentile of Download throughput [kbit/s]	27 500 kbit/s	16 500 kbit/s	88 000 kbit/s	66 000 kbit/s	8,0 %

B.4.2.3.2 File Upload (based on 1 MB File Size)

Factor	Lower limit		Upper limit		Weight
	Cities	Outside cities	Cities	Outside cities	
Composite Session Success Criterion combining: <ul style="list-style-type: none"> Success to access and download the file Achievement of min. throughput of 384 kbit/s 	87,0 %	80,0 %	100,0 %		60,0 %
Average Upload Session Duration [s]	0,5 s	1,0s	8,0s	10,0s	12,0 %
10 th percentile of Upload throughput [kbit/s]	384 kbit/s		8 800 kbit/s	6 600 kbit/s	20,0 %
90 th percentile of Upload throughput [kbit/s]	1 100,0 kbit/s	5 500 kbit/s	33 000 kbit/s	27 500 kbit/s	8,0 %

B.4.2.3.3 File Download (based on 7 s Fixed Download Time)

Factor	Lower limit		Upper limit		Weight
	Cities	Outside cities	Cities	Outside cities	
Composite Session Success Criterion combining: <ul style="list-style-type: none"> Success to access and download the file for the full duration Achievement of min. throughput of 384 kbit/s 	87,0 %	80,0 %	100,0 %		60,0 %
Average Download Throughput [kbit/s]	8 250 kbit/s	4 400 kbit/s	88 000 kbit/s	55 000 kbit/s	12,0 %
10 th percentile of Download throughput [kbit/s]	384 kbit/s		38 500 kbit/s	13 200 kbit/s	20,0 %
90 th percentile of Download throughput [kbit/s]	27 500 kbit/s	16 500 kbit/s	198 000 kbit/s	110 000 kbit/s	8,0 %

B.4.2.3.4 File Upload (based on 7 s Fixed Upload Time)

Factor	Lower limit		Upper limit		Weight
	Cities	Outside cities	Cities	Outside cities	
Composite Session Success Criterion combining: <ul style="list-style-type: none"> Success to access and upload the file Achievement of min. throughput of 384 kbit/s 	87,0 %	80,0 %	100,0 %		60,0 %
Average Upload Throughput [kbit/s]	4 400 kbit/s	1 100 kbit/s	50 000 kbit/s	35 000 kbit/s	12,0 %
10 th percentile of Upload throughput [kbit/s]	384 kbit/s		15 000 kbit/s	10 000 kbit/s	20,0 %
90 th percentile of Upload throughput [kbit/s]	16 500 kbit/s	6 600 kbit/s	75 000 kbit/s	53 000 kbit/s	8,0 %

B.4.2.4 Browsing

B.4.2.4.1 Static Web Page

Factor	Lower limit		Upper limit		Weight
	Cities	Outside cities	Cities	Outside cities	
Static Web Page (Kepler [i.3]) Composite Session Success: <ul style="list-style-type: none"> Success to access and upload the entire web page 	87,0 %	80,0 %	100,0 %		60,0 %
Average Session Duration [s]	1,0 s		4,0 s		40,0 %

B.4.2.4.2 Dynamic Web Pages

Factor	Lower limit		Upper limit		Weight
	Cities	Outside cities	Cities	Outside cities	
Dynamic Web Pages Composite Session Success: <ul style="list-style-type: none"> Success to access and upload the entire web page Achievement of min. throughput of 384 kbit/s in first second 	87,0 %	80,0 %	100,0 %		60,0 %
Combined DNS resolution and IP Service Access Duration	0,3 s		1,0 s	1,5 s	20,0 %
Throughput on first 800 kB	384 kbit/s		6 400 kbit/s		20,0 %

B.5 Remarks on mapping functions

The example presented in the above clauses assumes that different types of mapping functions are used for the different factors/QoS parameters.

Linear functions are used for all QoS parameters with the following exceptions, where square root-shaped mapping functions are used: Speech Quality, Data Throughput, Avg. Video Resolution, and File DL and UL Average Session Duration.

B.6 Example Calculation

	Common Parameters					Cities													Outside Cities		
	Unit	Weight in Data or Telephony	Weight in Data or Telephony	Weight in Service	Exponent (L=linear, 0.5=sqrt)	City/Town Limits		Big and Medium			Complementary Areas			towns		Road Limits		roads			
						Bad Limit	Good Limit	Example Result city drivetest	RAW = MIN(MAX((G-E)/(F-E))*100; 0);100)^D	Score = A*B*C* RAW	Example Result city complementary areas (walktest)	RAW = MIN(MAX((H-E)/(F-E))*100; 0);100)^D	Score = A*B*C* RAW	Example Result towns drivetest	RAW = MIN(MAX((I-E)/(F-E))*100; 0);100)^D	Score = A*B*C* RAW	Bad Limit	Good Limit	Example Result roads drivetest	RAW = MIN(MAX((L-E)/(K-E))*100; 0);100)^D	Score = A*B*C* RAW
A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T		
Telephony																					
COMPOSITE SUCCESS CRITERION	[%]	40%	100%	70%	1	87.0%	100.0%	98.2%	85.8%	24.0%	98.6%	89.0%	24.9%	98.1%	85.3%	23.9%	80.0%	100.0%	93.3%	66.5%	18.6%
AVG CALL SETUP TIME	[s]	40%	100%	3%	1	8	2.5	7.3	12.6%	0.2%	6.5	27.1%	0.3%	7.0	18.0%	0.2%	8	2.5	7.2	13.9%	0.2%
CALL SETUP TIME P90	[s]	40%	100%	12%	1	10	3	8.5	20.9%	1.0%	7.2	39.5%	1.9%	8.2	25.7%	1.2%	10	3	9.1	12.9%	0.6%
AVG SPEECH QUAL	[MOS]	40%	100%	15%	0.5	2.5	4.4	3.1	55.2%	3.3%	3.1	56.3%	3.4%	3.0	49.7%	3.0%	2.5	4.4	2.8	41.5%	2.5%
Data Services																					
Video Streaming																					
Youtube Video																					
COMPOSITE SUCCESS CRITERION	[%]	60%	22%	55%	1	87.0%	100.0%	99.3%	94.4%	6.9%	98.6%	89.4%	6.6%	99.1%	93.3%	6.8%	80.0%	100.0%	95.9%	79.3%	5.8%
VIDEO START TIME	[s]	60%	22%	5%	1	4.5	1.5	2.5	65.4%	0.4%	2.3	71.7%	0.5%	2.5	65.8%	0.4%	5	1.5	2.8	61.8%	0.4%
PLAYOUTS WO IN TERRUPTION S	[%]	60%	22%	15%	1	95.0%	100.0%	99.0%	79.9%	1.6%	100.0%	100.0%	2.0%	98.8%	76.6%	1.5%	90.0%	100.0%	95.8%	58.5%	1.2%
AVG RESOLUTION	[p]	60%	22%	25%	0.5	480	1080	1079	99.9%	3.3%	1063	98.5%	3.3%	1080	100.0%	3.3%	480	1080	1018	94.7%	3.2%
File Download																					
3MB Fixed Size																					
COMPOSITE SUCCESS CRITERION	[%]	60%	14%	60%	1	87.0%	100.0%	100.0%	100.0%	5.0%	99.8%	98.7%	4.9%	99.3%	94.8%	4.7%	80.0%	100.0%	97.8%	89.0%	4.4%
AVERAGE SESSION TIME	[s]	60%	14%	12%	0.5	10	1	2.2	93.3%	0.9%	2.4	92.2%	0.9%	2.2	93.0%	0.9%	10	1.5	3.7	86.3%	0.9%
P10 DATA RATE	[kbit/s]	60%	14%	20%	0.5	384	27500	6833	48.8%	0.8%	6685	48.2%	0.8%	6593	48.2%	0.8%	384	11000	3245	51.9%	0.9%
P90 DATA RATE	[kbit/s]	60%	14%	8%	0.5	27500	88000	36731	39.1%	0.3%	39474	44.5%	0.3%	35529	36.4%	0.2%	16500	66000	35134	61.4%	0.4%
7s Fixed Duration																					
COMPOSITE SUCCESS CRITERION	[%]	60%	14%	60%	1	87.0%	100.0%	99.9%	99.3%	5.0%	99.8%	98.7%	4.9%	99.9%	99.0%	4.9%	80.0%	100.0%	97.6%	87.9%	4.4%
AVERAGE DATA RATE	[kbit/s]	60%	14%	12%	0.5	8250	88000	26967	48.4%	0.5%	29932	52.1%	0.5%	28270	50.1%	0.5%	4400	55000	26440	66.0%	0.7%
P10 DATA RATE	[kbit/s]	60%	14%	20%	0.5	384	38500	8861	47.2%	0.8%	9666	49.3%	0.8%	10294	51.0%	0.8%	384	13200	5768	64.8%	1.1%
P90 DATA RATE	[kbit/s]	60%	14%	8%	0.5	27500	198000	47884	34.6%	0.2%	53867	39.3%	0.3%	48599	35.2%	0.2%	16500	110000	51643	61.3%	0.4%
File Upload																					
1MB Fixed Size																					
COMPOSITE SUCCESS CRITERION	[%]	60%	14%	60%	1	87.0%	100.0%	99.8%	98.6%	4.9%	98.8%	91.0%	4.5%	99.3%	94.7%	4.7%	80.0%	100.0%	96.3%	81.7%	4.1%
AVERAGE SESSION TIME	[s]	60%	14%	12%	0.5	8	0.5	1.1	95.7%	1.0%	1.3	94.6%	0.9%	1.4	94.1%	0.9%	10	1	1.8	95.3%	1.0%
P10 DATA RATE	[kbit/s]	60%	14%	20%	0.5	384	8800	4817	72.6%	1.2%	3537	61.2%	1.0%	3581	61.6%	1.0%	384	6600	2275	55.2%	0.9%
P90 DATA RATE	[kbit/s]	60%	14%	8%	0.5	11000	33000	20202	64.7%	0.4%	21517	69.1%	0.5%	18349	57.8%	0.4%	5500	27500	17529	73.9%	0.5%
7s Fixed Duration																					
COMPOSITE SUCCESS CRITERION	[%]	60%	14%	60%	1	87.0%	100.0%	99.6%	96.6%	4.8%	98.5%	88.5%	4.4%	99.0%	92.7%	4.6%	80.0%	100.0%	95.1%	80.4%	4.0%
AVERAGE DATA RATE	[kbit/s]	60%	14%	12%	0.5	4400	50000	20567	59.5%	0.6%	20567	59.5%	0.6%	18291	55.2%	0.6%	1100	35000	16104	66.5%	0.7%
P10 DATA RATE	[kbit/s]	60%	14%	20%	0.5	384	15000	5883	61.3%	1.0%	4350	52.1%	0.9%	3581	46.8%	0.8%	384	10000	2555	47.5%	0.8%
P90 DATA RATE	[kbit/s]	60%	14%	8%	0.5	16500	75000	35402	56.8%	0.4%	38658	61.5%	0.4%	35362	56.8%	0.4%	6600	53000	31366	73.1%	0.5%
Web Browsing																					
Kepler Static Page																					
COMPOSITE SUCCESS CRITERION	[%]	60%	0%	60%	1	87.0%	100.0%	99.8%	98.8%	2.0%	99.6%	96.8%	1.9%	99.5%	96.4%	1.9%	80.0%	100.0%	96.6%	83.2%	1.7%
OVERALL SESSION TIME	[s]	60%	0%	40%	1	4	1	1.3	90.6%	1.2%	1.2	92.2%	1.2%	1.3	89.0%	1.2%	4	1	1.7	76.7%	1.0%
Live Web Pages																					
COMPOSITE SUCCESS CRITERION	[%]	60%	17%	60%	1	87.0%	100.0%	99.5%	96.3%	5.8%	99.2%	93.9%	5.6%	99.3%	94.8%	5.7%	80.0%	100.0%	96.6%	82.8%	5.0%
OVERALL SESSION TIME	[s]	60%	17%	40%	1	4	1	2.0	67.2%	2.7%	2.0	65.5%	2.6%	2.3	55.7%	2.2%	4	1	2.7	43.4%	1.7%
									Sum (City)	80.3%	Sum (Complementary)		81.0%	Sum (Towns)		78.2%	Sum (Roads)			67.4%	
									Weight	45%	Weight		15%	Weight		20%	Weight			20%	
																			Total Score	77.4%	

History

Document history		
V1.1.1	August 2019	Publication