

ETSI TS 128 105 V18.9.0 (2026-04)



TECHNICAL SPECIFICATION

**5G;
Management and orchestration;
Artificial Intelligence/ Machine Learning (AI/ML) management
(3GPP TS 28.105 version 18.9.0 Release 18)**

get full document from standards.iteh.ai



Reference

RTS/TSGS-0528105vi90

Keywords

5G

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 - APE 7112B
Association à but non lucratif enregistrée à la
Sous-Préfecture de Grasse (06) N° w061004871

Important notice

The present document can be downloaded from the
[ETSI Search & Browse Standards](#) application.

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 on [ETSI deliver](#) repository.

Users should be aware that the present document may be revised or have its status changed, this information is available in the [Milestones listing](#).

If you find errors in the present document, please send your comments to the relevant service listed under [Committee Support Staff](#).

If you find a security vulnerability in the present document, please report it through our [Coordinated Vulnerability Disclosure \(CVD\)](#) program.

Notice of disclaimer & limitation of liability

The information provided in the present deliverable is directed solely to professionals who have the appropriate degree of experience to understand and interpret its content in accordance with generally accepted engineering or other professional standard and applicable regulations.

No recommendation as to products and services or vendors is made or should be implied.

No representation or warranty is made that this deliverable is technically accurate or sufficient or conforms to any law and/or governmental rule and/or regulation and further, no representation or warranty is made of merchantability or fitness for any particular purpose or against infringement of intellectual property rights.

In no event shall ETSI be held liable for loss of profits or any other incidental or consequential damages.

Any software contained in this deliverable is provided "AS IS" with no warranties, express or implied, including but not limited to, the warranties of merchantability, fitness for a particular purpose and non-infringement of intellectual property rights and ETSI shall not be held liable in any event for any damages whatsoever (including, without limitation, damages for loss of profits, business interruption, loss of information, or any other pecuniary loss) arising out of or related to the use of or inability to use the software.

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 2026.
All rights reserved.

Intellectual Property Rights

Essential patents

IPRs essential or potentially essential to normative deliverables may have been declared to ETSI. The declarations pertaining to these essential IPRs, if any, are 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 IPR online database](#).

Pursuant to the ETSI Directives including the ETSI IPR Policy, no investigation regarding the essentiality of IPRs, 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.

DECT™, **PLUGTESTS™**, **UMTS™** and the ETSI logo are trademarks of ETSI registered for the benefit of its Members. **3GPP™**, **LTE™** and **5G™** logo 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.

Legal Notice

This Technical Specification (TS) has been produced by ETSI 3rd Generation Partnership Project (3GPP).

The present document may refer to technical specifications or reports using their 3GPP identities. These shall be interpreted as being references to the corresponding ETSI deliverables.

The cross reference between 3GPP and ETSI identities can be found at [3GPP to ETSI numbering cross-referencing](#).

Modal verbs terminology

In the present document "**shall**", "**shall not**", "**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.

Contents

Intellectual Property Rights	2
Legal Notice	2
Modal verbs terminology.....	2
Foreword.....	9
1 Scope	11
2 References	11
3 Definitions of terms, symbols and abbreviations	12
3.1 Terms.....	12
3.2 Symbols.....	13
3.3 Abbreviations	13
4 Concepts and overview	13
4.1 Overview	13
4a AI/ML management functionality and service framework	14
4a.0 ML model lifecycle	14
4a.1 Functionality and service framework for ML model training	15
4a.2 AI/ML functionalities management scenarios (relation with managed AI/ML features).....	16
5 Void.....	18
6 AI/ML management use cases and requirements	18
6.1 ML model lifecycle management capabilities	18
6.2 Void.....	19
6.2a Void.....	19
6.2b ML model training.....	19
6.2b.1 Description.....	19
6.2b.2 Use cases.....	20
6.2b.2.1 ML model training requested by consumer.....	20
6.2b.2.2 ML model training initiated by producer	20
6.2b.2.3 ML model selection.....	21
6.2b.2.4 Managing ML model training processes	21
6.2b.2.5 Handling errors in data and ML decisions	21
6.2b.2.6 ML model joint training	22
6.2b.2.7 ML model validation performance reporting	22
6.2b.2.8 Training data effectiveness reporting.....	23
6.2b.2.9 Performance management for ML model training	23
6.2b.2.9.1 Overview	23
6.2b.2.9.2 Performance indicator selection for ML model training.....	23
6.2b.2.9.3 ML model performance indicators query and selection for ML model training.....	23
6.2b.2.9.4 MnS consumer policy-based selection of ML model performance indicators for ML model training.....	24
6.2b.3 Requirements for ML model training	24
6.2c ML model testing	26
6.2c.1 Description.....	26
6.2c.2 Use cases.....	27
6.2c.2.1 Consumer-requested ML model testing	27
6.2c.2.2 Producer-initiated ML model testing	27
6.2c.2.3 Joint testing of multiple ML models	27
6.2c.2.4 Performance management for ML model testing	27
6.2c.2.4.1 Overview	27
6.2c.2.4.2 Performance indicator selection for ML model testing	27
6.2c.2.4.3 ML model performance indicators query and selection for ML model testing	28
6.2c.2.4.4 MnS consumer policy-based selection of ML model performance indicators for ML model testing	28
6.2c.3 Requirements for ML model testing	28

6.3	AI/ML inference emulation	29
6.3.1	Description	29
6.3.2	Use cases	29
6.3.2.1	AI/ML inference emulation	29
6.3.3	Requirements for Managing AI/ML inference emulation	29
6.4	ML model deployment	29
6.4.1	ML model loading	29
6.4.1.1	Description	29
6.4.1.2	Use cases	30
6.4.1.2.1	Consumer requested ML model loading	30
6.4.1.2.2	Control of producer-initiated ML model loading	30
6.4.1.2.3	ML model registration	30
6.4.1.3	Requirements for ML model loading	30
6.5	AI/ML inference	31
6.5.1	AI/ML inference performance management	31
6.5.1.1	Description	31
6.5.1.2	Use cases	31
6.5.1.2.1	AI/ML inference performance evaluation	31
6.5.1.2.2	AI/ML performance measurements selection based on MnS consumer policy	32
6.5.1.3	Requirements for AI/ML inference performance management	32
6.5.2	AI/ML update control	32
6.5.2.1	Description	32
6.5.2.2	Use cases	33
6.5.2.2.1	Availability of new capabilities or ML models	33
6.5.2.2.2	Triggering ML model update	33
6.5.2.3	Requirements for AI/ML update control	33
6.5.3	AI/ML inference capabilities management	34
6.5.3.1	Description	34
6.5.3.2	Use cases	34
6.5.3.2.1	Identifying capabilities of ML models	34
6.5.3.2.2	Mapping of the capabilities of ML models	35
6.5.3.3	Requirements for AI/ML inference capabilities management	35
6.5.4	AI/ML inference capability configuration management	35
6.5.4.1	Description	35
6.5.4.2	Use cases	35
6.5.4.2.1	Managing NG-RAN AI/ML-based distributed Network Energy Saving	35
6.5.4.2.2	Managing NG-RAN AI/ML-based distributed Mobility Optimization	36
6.5.4.2.3	Managing NG-RAN AI/ML-based distributed Load Balancing	36
6.5.4.3	Requirements for AI/ML inference management	36
6.5.5	AI/ML Inference History	37
6.5.5.1	Description	37
6.5.5.2	Use cases	37
6.5.5.2.1	AI/ML Inference History - tracking inferences and context	37
6.5.5.3	Requirements for AI/ML Inference History	37
7	Information model definitions for AI/ML management	38
7.1	Imported and associated information entities	38
7.1.1	Imported information entities and local labels	38
7.1.2	Associated information entities and local labels	38
7.2	Void	38
7.2a	Common information model definitions for AI/ML management	38
7.2a.1	Class diagram	38
7.2a.1.1	Relationships	38
7.2a.1.2	Inheritance	39
7.2a.2	Class definitions	39
7.2a.2.1	MLModel	39
7.2a.2.1.1	Definition	39
7.2a.2.1.2	Attributes	40
7.2a.2.1.3	Attribute constraints	40
7.2a.2.1.4	Notifications	40
7.2a.2.2	MLModelRepository	40
7.2a.2.2.1	Definition	40

7.2a.2.2.2	Attributes	40
7.2a.2.2.3	Attribute constraints	40
7.2a.2.2.4	Notifications	40
7.2a.2.3	MLModelCoordinationGroup	40
7.2a.2.3.1	Definition.....	40
7.2a.2.3.2	Attributes.....	41
7.2a.2.3.3	Attribute constraints	41
7.2a.2.3.4	Notifications	41
7.3	Void.....	41
7.3a	Information model definitions for AI/ML operational phases.....	41
7.3a.1	Information model definitions for ML model training	41
7.3a.1.1	Class diagram.....	41
7.3a.1.1.1	Relationships	41
7.3a.1.1.2	Inheritance	42
7.3a.1.2	Class definitions.....	42
7.3a.1.2.1	MLTrainingFunction.....	42
7.3a.1.2.1.1	Definition.....	42
7.3a.1.2.1.2	Attributes	43
7.3a.1.2.1.3	Attribute constraints.....	43
7.3a.1.2.1.4	Notifications.....	43
7.3a.1.2.2	MLTrainingRequest	43
7.3a.1.2.2.1	Definition.....	43
7.3a.1.2.2.2	Attributes	44
7.3a.1.2.2.3	Attribute constraints.....	44
7.3a.1.2.2.4	Notifications.....	44
7.3a.1.2.3	MLTrainingReport.....	44
7.3a.1.2.3.1	Definition.....	44
7.3a.1.2.3.2	Attributes	45
7.3a.1.2.3.3	Attribute constraints.....	45
7.3a.1.2.3.4	Notifications.....	45
7.3a.1.2.4	MLTrainingProcess	45
7.3a.1.2.4.1	Definition.....	45
7.3a.1.2.4.2	Attributes	46
7.3a.1.2.4.3	Attribute constraints.....	46
7.3a.1.2.4.4	Notifications.....	47
7.3a.2	Information model definitions for AI/ML inference emulation.....	50
7.3a.2.1	Class diagram.....	50
7.3a.2.1.1	Relationships	50
7.3a.2.1.2	Inheritance	50
7.3a.2.2	Class definitions.....	51
7.3a.2.2.1	AIMLInferenceEmulationFunction	51
7.3a.2.2.1.1	Definition.....	51
7.3a.2.2.1.2	Attributes	51
7.3a.2.2.1.3	Attribute constraints.....	51
7.3a.2.2.1.4	Notifications.....	51
7.3a.3	Information model definitions for ML model deployment	51
7.3a.3.1	Class diagram.....	51
7.3a.3.1.1	Relationships	51
7.3a.3.1.2	Inheritance	52
7.3a.3.2	Class definitions.....	52
7.3a.3.2.1	MLModelLoadingRequest.....	52
7.3a.3.2.1.1	Definition.....	52
7.3a.3.2.1.2	Attributes	53
7.3a.3.2.1.3	Attribute constraints.....	53
7.3a.3.2.1.4	Notifications.....	53
7.3a.3.2.2	MLModelLoadingPolicy	53
7.3a.3.2.2.1	Definition.....	53
7.3a.3.2.2.2	Attributes	53
7.3a.3.2.2.3	Attribute constraints.....	53
7.3a.3.2.2.4	Notifications.....	53
7.3a.3.2.3	MLModelLoadingProcess.....	54

7.3a.3.2.3.1	Definition	54
7.3a.3.2.3.2	Attributes	54
7.3a.3.2.3.3	Attribute constraints	55
7.3a.3.2.3.4	Notifications	55
7.3a.4	Information model definitions for ML inference	55
7.3a.4.1	Class diagram	55
7.3a.4.1.1	Relationships	55
7.3a.4.1.2	Inheritance	56
7.3a.4.2	Class definitions	56
7.3a.4.2.1	MLUpdateFunction	56
7.3a.4.2.1.1	Definition	56
7.3a.4.2.1.2	Attributes	57
7.3a.4.2.1.3	Attribute constraints	57
7.3a.4.2.1.4	Notifications	57
7.3a.4.2.2	MLUpdateRequest	57
7.3a.4.2.2.1	Definition	57
7.3a.4.2.2.2	Attributes	58
7.3a.4.2.2.3	Attribute constraints	58
7.3a.4.2.2.4	Notifications	58
7.3a.4.2.3	MLUpdateProcess	58
7.3a.4.2.3.1	Definition	58
7.3a.4.2.3.2	Attributes	59
7.3a.4.2.3.3	Attribute constraints	59
7.3a.4.2.3.4	Notifications	59
7.3a.4.2.4	MLUpdateReport	59
7.3a.4.2.4.1	Definition	59
7.3a.4.2.4.2	Attributes	60
7.3a.4.2.4.3	Attribute constraints	60
7.3a.4.2.4.4	Notifications	60
7.3a.4.2.5	AIMLInferenceFunction	60
7.3a.4.2.5.1	Definition	60
7.3a.4.2.5.2	Attributes	61
7.3a.4.2.5.3	Attribute constraints	61
7.3a.4.2.5.4	Notifications	61
7.3a.4.2.6	AIMLInferenceReport	61
7.3a.4.2.6.1	Definition	61
7.3a.4.2.6.2	Attributes	61
7.3a.4.2.6.3	Attribute constraints	61
7.3a.4.2.6.4	Notifications	61
7.4	Data type definitions	62
7.4.1	ModelPerformance <<dataType>>	62
7.4.1.1	Definition	62
7.4.1.2	Attributes	62
7.4.1.3	Attribute constraints	62
7.4.1.4	Notifications	62
7.4.2	Void	62
7.4.3	MLContext <<dataType>>	62
7.4.3.1	Definition	62
7.4.3.2	Attributes	62
7.4.3.3	Attribute constraints	62
7.4.3.4	Notifications	63
7.4.4	SupportedPerfIndicator <<dataType>>	63
7.4.4.1	Definition	63
7.4.4.2	Attributes	63
7.4.4.3	Attribute constraints	63
7.4.4.4	Notifications	63
7.4.5	AvailMLCapabilityReport <<dataType>>	63
7.4.5.1	Definition	63
7.4.5.2	Attributes	64
7.4.5.3	Attribute constraints	64
7.4.5.4	Notifications	64

7.4.6	AIMLManagementPolicy <<dataType>>.....	64
7.4.6.1	Definition	64
7.4.6.2	Attributes.....	64
7.4.6.3	Attribute constraints	64
7.4.6.4	Notifications.....	64
7.4.7	ManagedActivationScope <<choice>>	64
7.4.7.1	Definition	64
7.4.7.2	Attributes.....	65
7.4.7.3	Attribute constraints	65
7.4.7.4	Notifications.....	65
7.4.8	MLCapabilityInfo <<dataType>>	65
7.4.8.1	Definition	65
7.4.8.2	Attributes.....	65
7.4.8.3	Attribute constraints	65
7.4.8.4	Notifications.....	65
7.4.9	InferenceOutput <<dataType>>	65
7.4.9.1	Definition	65
7.4.9.2	Attributes.....	66
7.4.9.3	Attribute constraints	66
7.4.9.4	Notifications.....	66
7.4.10	AIMLInferenceName <<choice>>	66
7.4.10.1	Definition	66
7.4.10.2	Attributes.....	66
7.4.10.3	Attribute constraints	66
7.4.10.4	Notifications.....	66
7.4a	Enumerations.....	67
7.4a.1	NgRanInferenceType <<enumeration>>.....	67
7.5	Attribute definitions	68
7.5.1	Attribute properties	68
7.5.2	Constraints	81
7.6	Common notifications	81
7.6.1	Configuration notifications	81
8	Service components.....	81
8.0	General	81
8.1	Lifecycle management operations for AI/ML management MnS	81
9	Solution Set (SS)	83
Annex A (informative): PlantUML source code for NRM class diagrams.....		84
A.1	General	84
A.2	PlantUML code for Figure 7.3a.1.1.1-1: NRM fragment for ML model training.....	84
A.3	PlantUML code for Figure 7.3a.1.1.2-1: Inheritance Hierarchy for ML model training related NRMs	85
A.4	PlantUML code for Figure 7.2a.1.2-1: Inheritance Hierarchy for common information models for AI/ML management	86
A.5	PlantUML code for Figure 7.2a.1.1-1: Relationships for common information models for AI/ML management	86
A.6	PlantUML code for Figure 7.3a.1.1.1-2: NRM fragment for ML model testing	86
A.7	PlantUML code for Figure 7.3a.1.1.2-2: Inheritance Hierarchy for ML model testing related NRMs	87
A.8	PlantUML code for Figure 7.3a.4.1.1-1: NRM fragment for ML update	87
A.9	PlantUML code for Figure 7.3a.4.1.2-1: Inheritance Hierarchy for ML update related NRMs.....	88
A.10	PlantUML code for Figure 7.3a.3.1.1-1: NRM fragment for ML model loading	88

A.11 PlantUML code for Figure 7.3a.3.1.2-1: Inheritance Hierarchy for ML model loading related NRMs	89
A.12 PlantUML code for Figure 7.3a.4.1.1-2: NRM fragment for AI/ML inference function.....	89
A.13 PlantUML code for Figure 7.3a.4.1.2-2: Inheritance Hierarchy for AI/ML inference function	90
A.14 PlantUML code for Figure 7.3a.2.1.1-1: NRM fragment for AI/ML inference emulation Control.....	90
A.15 PlantUML code for Figure 7.3a.2.1.2-1: AI/ML inference emulation Inheritance Relations	91
Annex B (normative): OpenAPI definition of the AI/ML NRM	92
B.1 General	92
B.2 Solution Set (SS) definitions	92
B.2.1 OpenAPI document "TS28105_AiMLNrm.yaml"	92
Annex C (informative): Change history	93
History	96

Sample Document

get full document from standards.iteh.ai

Foreword

This Technical Specification has been produced by the 3rd Generation Partnership Project (3GPP).

The contents of the present document are subject to continuing work within the TSG and may change following formal TSG approval. Should the TSG modify the contents of the present document, it will be re-released by the TSG with an identifying change of release date and an increase in version number as follows:

Version x.y.z

where:

- x the first digit:
 - 1 presented to TSG for information;
 - 2 presented to TSG for approval;
 - 3 or greater indicates TSG approved document under change control.
- y the second digit is incremented for all changes of substance, i.e. technical enhancements, corrections, updates, etc.
- z the third digit is incremented when editorial only changes have been incorporated in the document.

In the present document, modal verbs have the following meanings:

shall indicates a mandatory requirement to do something

shall not indicates an interdiction (prohibition) to do something

The constructions "shall" and "shall not" are confined to the context of normative provisions, and do not appear in Technical Reports.

The constructions "must" and "must not" are not used as substitutes for "shall" and "shall not". Their use is avoided insofar as possible, and they are not used in a normative context except in a direct citation from an external, referenced, non-3GPP document, or so as to maintain continuity of style when extending or modifying the provisions of such a referenced document.

should indicates a recommendation to do something

should not indicates a recommendation not to do something

may indicates permission to do something

need not indicates permission not to do something

The construction "may not" is ambiguous and is not used in normative elements. The unambiguous constructions "might not" or "shall not" are used instead, depending upon the meaning intended.

can indicates that something is possible

cannot indicates that something is impossible

The constructions "can" and "cannot" are not substitutes for "may" and "need not".

will indicates that something is certain or expected to happen as a result of action taken by an agency the behaviour of which is outside the scope of the present document

will not indicates that something is certain or expected not to happen as a result of action taken by an agency the behaviour of which is outside the scope of the present document

might indicates a likelihood that something will happen as a result of action taken by some agency the behaviour of which is outside the scope of the present document

might not indicates a likelihood that something will not happen as a result of action taken by some agency the behaviour of which is outside the scope of the present document

In addition:

is (or any other verb in the indicative mood) indicates a statement of fact

is not (or any other negative verb in the indicative mood) indicates a statement of fact

The constructions "is" and "is not" do not indicate requirements.

Sample Document

get full document from standards.iteh.ai

1 Scope

The present document specifies the Artificial Intelligence / Machine Learning (AI/ML) management capabilities and services for 5GS where AI/ML is used, including management and orchestration (e.g., MDA, see 3GPP TS 28.104 [2]) and 5G networks (e.g. NWDAF, see 3GPP TS 23.288 [3]) and NG-RAN (see TS 38.300 [16] and TS 38.401 [17]).

2 References

The following documents contain provisions which, through reference in this text, constitute provisions of the present document.

- References are either specific (identified by date of publication, edition number, version number, etc.) or non-specific.
- For a specific reference, subsequent revisions do not apply.
- For a non-specific reference, the latest version applies. In the case of a reference to a 3GPP document (including a GSM document), a non-specific reference implicitly refers to the latest version of that document *in the same Release as the present document*.

- [1] 3GPP TR 21.905: "Vocabulary for 3GPP Specifications".
- [2] 3GPP TS 28.104: "Management and orchestration; Management Data Analytics".
- [3] 3GPP TS 23.288: "Architecture enhancements for 5G System (5GS) to support network data analytics services".
- [4] 3GPP TS 28.552: "Management and orchestration; 5G performance measurements".
- [5] 3GPP TS 32.425: "Telecommunication management; Performance Management (PM); Performance measurements Evolved Universal Terrestrial Radio Access Network (E-UTRAN)".
- [6] 3GPP TS 28.554: "Management and orchestration; 5G end to end Key Performance Indicators (KPI)".
- [7] 3GPP TS 32.422: "Telecommunication management; Subscriber and equipment trace; Trace control and configuration management".
- [8] Void
- [9] 3GPP TS 28.405: "Telecommunication management; Quality of Experience (QoE) measurement collection; Control and configuration".
- [10] Void
- [11] 3GPP TS 28.532: "Management and orchestration; Generic management services".
- [12] 3GPP TS 28.622: "Telecommunication management; Generic Network Resource Model (NRM) Integration Reference Point (IRP); Information Service (IS)".
- [13] 3GPP TS 32.156: "Telecommunication management; Fixed Mobile Convergence (FMC) Model repertoire".
- [14] 3GPP TS 32.160: "Management and orchestration; Management service template".
- [15] 3GPP TS 28.533: "Management and orchestration; Architecture framework".
- [16] 3GPP TS 38.300: "NR; NR and NG-RAN Overall description; Stage-2".
- [17] 3GPP TS 38.401: "NG-RAN; Architecture description".
- [18] 3GPP TS 28.541: " Management and orchestration; 5G Network Resource Model (NRM); Stage 2 and stage 3".

- [19] 3GPP TS 28.623: "Telecommunication management; Generic Network Resource Model (NRM) Integration Reference Point (IRP); Solution Set (SS) definitions".
- [20] 3GPP TS 29.520: "5G System; Network Data Analytics Services; Stage 3".

3 Definitions of terms, symbols and abbreviations

3.1 Terms

For the purposes of the present document, the terms given in 3GPP TR 21.905 [1] and the following apply. A term defined in the present document takes precedence over the definition of the same term, if any, in 3GPP TR 21.905 [1].

ML model: a manageable representation of an ML model algorithm.

NOTE 1: an ML model algorithm is a mathematical algorithm through which running a set of input data can generate a set of inference output.

NOTE 2: ML model algorithm is proprietary and not in scope for standardization and therefore not treated in this specification.

NOTE 3: ML model may include metadata. Metadata may include e.g. information related to the trained model, and applicable runtime context.

ML model training: a process performed by an ML training function to take training data, run it through an ML model algorithm, derive the associated loss and adjust the parameterization of that ML model iteratively based on the computed loss and generate the trained ML model.

ML model initial training: a process of training an initial version of an ML model.

ML model re-training: a process of training a previously trained version of an ML model and generate a new version.

NOTE 4: a new version of a trained ML model supports the same type of inference as the previous version of the ML model, i.e., the data type of inference input and data type of inference output remain unchanged between the two versions of the ML model, but parameter values might be different for the re-trained model.

ML model joint training: a process of training a group of ML models.

ML training function: a logical function with ML model training capabilities.

ML model testing: a process of evaluating the performance of an ML model using testing data different from data used for model training and validation.

ML model joint testing: a process of evaluating the performance of a group of ML models using testing data different from data used for model training and validation.

ML testing function: a logical function with ML model testing capabilities.

AI/ML inference: a process of running a set of input data through a trained ML model to produce set of output data, such as predictions.

NOTE 5: the inference represents the process to realize the AI capabilities by utilizing a trained ML model and other AI enablers if needed, hence the AI/ML prefix is used when referring to inference as compared to training and testing.

AI/ML inference function: a logical function that employs trained ML model(s) to conduct inference.

AI/ML inference emulation: running the inference process to evaluate the performance of an ML model in an emulation environment before deploying it into the target environment.

ML model deployment: a process of making a trained ML model available for use in the target environment.

ML model loading: a process of making a trained ML model available to an inference function.

AI/ML activation: a process of enabling the inference capability of an AI/ML inference function.

AI/ML deactivation: a process of disabling the inference capability of an AI/ML inference function.

3.2 Symbols

Void.

3.3 Abbreviations

For the purposes of the present document, the abbreviations given in TR 21.905 [1] and TS 28.533 [15]. An abbreviation defined in the present document takes precedence over the definition of the same abbreviation, if any, in TR 21.905 [1] and TS 28.533 [15].

AI	Artificial Intelligence
ML	Machine Learning

4 Concepts and overview

4.1 Overview

The AI/ML techniques and relevant applications are being increasingly adopted by the wider industries and proved to be successful. These are now being applied to telecommunication industry including mobile networks.

Although AI/ML techniques in general are quite mature nowadays, some of the relevant aspects of the technology are still evolving while new complementary techniques are frequently emerging.

The AI/ML techniques can be generally characterized from different perspectives including the followings:

- **Learning methods**

The learning methods include supervised learning, semi-supervised learning, unsupervised learning and reinforcement learning. Each learning method fits one or more specific category of inference (e.g. prediction), and requires specific type of training data. A brief comparison of these learning methods is provided in table 4.1-1.

Table 4.1-1: Comparison of Learning methods

	Supervised learning	Semi-supervised learning	Unsupervised learning	Reinforcement learning
Category of inference	Regression (numeric), classification	Regression (numeric), classification	Association, Clustering	Reward-based behaviour
Type of training data	Labelled data (Note)	Labelled data (Note), and unlabelled data	Unlabelled data	Not pre-defined
NOTE:	The labelled data refers to a set of training and testing data that have been assigned with one or more labels in order to add context and meaning.			

- **Learning complexity:**

- As per the learning complexity, there are Machine Learning (i.e. basic learning) and Deep Learning.

- **Learning architecture**

- Based on the topology and location where the learning tasks take place, the AI/ML can be categorized to centralized learning, distributed learning and federated learning.

- **Learning continuity**

- From learning continuity perspective, the AI/ML can be offline learning or continual learning.

Artificial Intelligence/Machine Learning (AI/ML) capabilities are used in various domains in 5GS, including management and orchestration (e.g. MDA, see 3GPP TS 28.104 [2]) and 5G networks (e.g. NWDAF, see 3GPP TS 23.288 [3]).

The AI/ML inference function in the 5GS uses the ML model for inference.

Each AI/ML technique, depending on the adopted specific characteristics as mentioned above, may be suitable for supporting certain type/category of use case(s) in 5GS.

To enable and facilitate the AI/ML capabilities with the suitable AI/ML techniques in 5GS, the ML model and AI/ML inference function need to be managed.

The present document specifies the generic AI/ML management related capabilities and services without specifically taking any of the above-mentioned learning methods into consideration. The AI/ML management capabilities which include the followings:

- ML model training.
- ML model testing.
- AI/ML inference emulation.
- ML model deployment.
- AI/ML inference.

4a AI/ML management functionality and service framework

4a.0 ML model lifecycle

AI/ML techniques are widely used in 5GS (including 5GC, NG-RAN, and management system), the generic AI/ML operational workflow shown in Figure 4a.0-1, highlights the main steps of an ML model lifecycle.

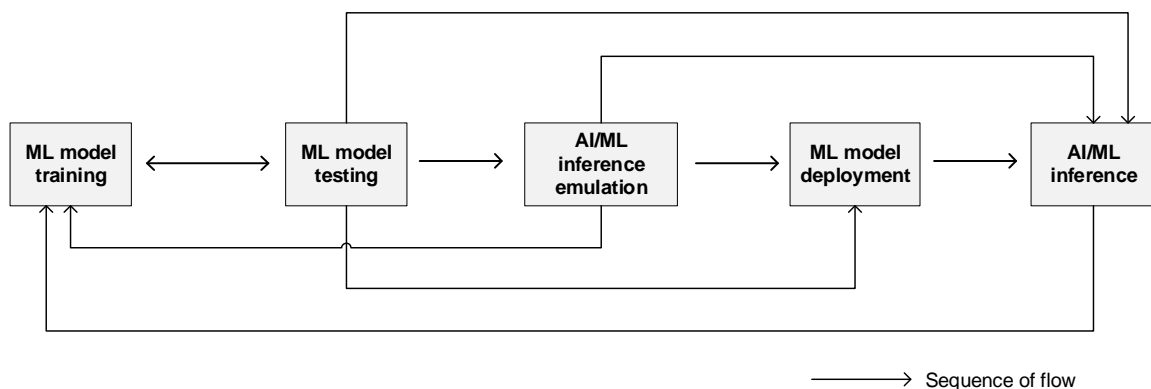


Figure 4a.0-1: ML model lifecycle

The ML model lifecycle includes training, testing, emulation, deployment, and inference. These steps are briefly described below:

- **ML model training:** training, including initial training and re-training, of an ML model or a group of ML models. It also includes validation of the ML model to evaluate the performance when the ML model performs on the training