



Moving Picture, Audio and Data Coding by  
Artificial Intelligence  
[www.mpai.community](http://www.mpai.community)

**N1001** 2022/12/21  
**Source** NNW-DC  
**Title** Technical Specification – Neural Network Watermarking (MPAI-NNW) WD 0.3  
**Target** MPAI Members

This document is a draft Neural Network Watermarking Technical Specification published for Community Comments.

Anybody can send comments to the [MPAI Secretariat](#) about this WD until 2023/01/23T15 UTC. MPAI plans on approving and publishing MMM Version 1 as a Technical Report on 2023/01/25.



Moving Picture, Audio and Data Coding  
by Artificial Intelligence  
[www.mpai.community](http://www.mpai.community)

## **MPAI Technical Specification**

### **Neural Network Watermarking MPAI-NNW**

**WD 0.3**

#### **WARNING**

Use of the technologies described in this Technical Specification may infringe patents, copyrights or intellectual property rights of MPAI Members or non-members.

MPAI and its Members accept no responsibility whatsoever for damages or liability, direct or consequential, which may result from the use of this Technical Specification.

Readers are invited to review Annex 2 – Notices and Disclaimers.

# Neural Network watermarking

## V1 [under development]

1	Introduction (Informative).....	3
2	Scope of Standard.....	5
3	Terms and Definitions.....	5
4	Use cases (Informative).....	6
5	References.....	6
5.1	Normative references.....	6
5.2	Informative references.....	6
6	Imperceptibility evaluation.....	6
7	Robustness evaluation.....	7
8	Computational cost evaluation.....	9
9	Example of usage (informative).....	<b>Error! Bookmark not defined.</b>
Annex 1	MPAI-wide terms and definitions.....	11
Annex 2	Notices and Disclaimers Concerning MPAI Standards (Informative).....	14
Annex 3	The Governance of the MPAI Ecosystem (Informative).....	16
Annex 4	Patent declarations.....	18
Annex 5	Imperceptibility evaluation.....	19
9.1.1	Classification task.....	19
9.1.2	Image/speech processing tasks.....	19
9.1.3	Image semantic segmentation.....	19

### 1 Introduction (Informative)

In recent years, Artificial Intelligence (AI) and related technologies have been introduced in a broad range of applications, have started affecting the life of millions of people and are expected to do so even more in the future. As digital media standards have positively influenced industry and billions of people, so AI-based data coding standards are expected to have a similar positive impact. Indeed, research has shown that data coding with AI-based technologies is generally *more efficient* than with existing technologies for, e.g., compression and feature-based description.

However, some AI technologies may carry inherent risks, e.g., in terms of bias toward some classes of users. Therefore, the need for standardisation is more important and urgent than ever.

The international, unaffiliated, not-for-profit MPAI – Moving Picture, Audio and Data Coding by Artificial Intelligence Standards Developing Organisation has the mission to develop *AI-enabled data coding standards*. MPAI Application Standards enable the development of AI-based products, applications, and services.

As a rule, MPAI standards include four documents: Technical Specification, Reference Software Specifications, Conformance Testing Specifications, and Performance Assessment Specifications. Sometimes Technical Reports are produced to provide informative guidance in specific areas for which the development of standards is premature.

Performance Assessment Specifications include standard operating procedures to enable users of MPAI Implementations to make informed decision about their applicability based on the notion of Performance, defined as a set of attributes characterising a reliable and trustworthy implementation. In the following, Terms beginning with a capital letter are defined in *Table 1* if they are specific to this Standard and in *Table 4* if they are common to all MPAI Standards.

In general, MPAI Application Standards are defined as aggregations – called AI Workflows (AIW) – of processing elements – called AI Modules (AIM) – executed in an AI Framework (AIF). MPAI

defines Interoperability as the ability to replace an AIW or an AIM Implementation with a functionally equivalent Implementation.

MPAI also defines 3 Interoperability Levels of an AIF that executes an AIW. The AIW and its AIMS may have 3 Levels:

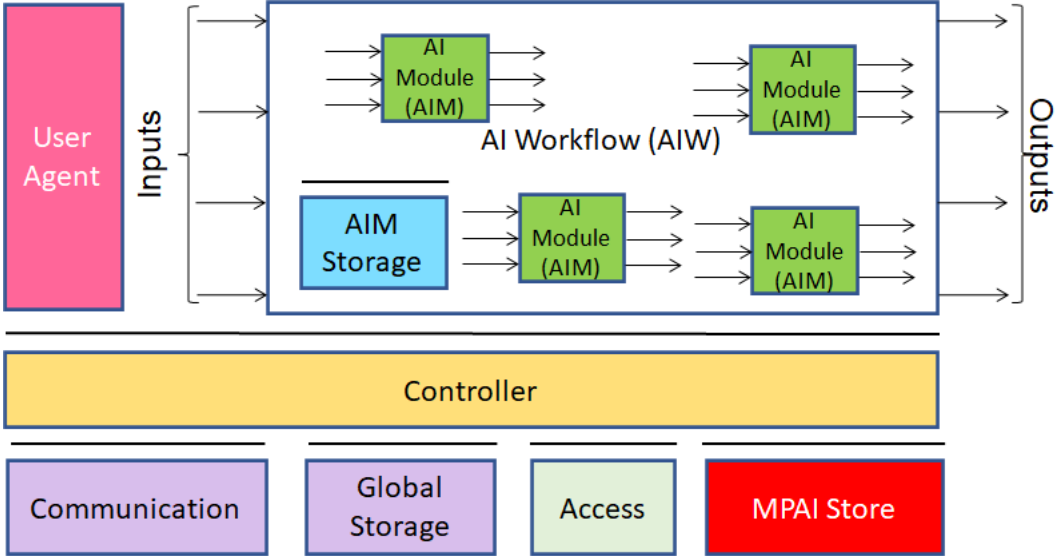
*Level 1* – Implementer-specific and satisfying the MPAI-AIF Standard.

*Level 2* – Specified by an MPAI Application Standard.

*Level 3* – Specified by an MPAI Application Standard and certified by a Performance Assessor.

MPAI offers Users access to the promised benefits of AI with a guarantee of increased transparency, trust and reliability as the Interoperability Level of an Implementation moves from 1 to 3. Additional information on Interoperability Levels is provided in reference [4].

*Figure 1* depicts the MPAI-AIF Reference Model under which Implementations of MPAI Application Standards and user-defined MPAI-AIF Conforming applications operate. MPAI is currently developing MPAI-AIF V2 that will compatibly extend MPAI-AIF V1 with security support.



*Figure 1 – The AI Framework (AIF) Reference Model and its Components*

MPAI Application Standards normatively specify the Syntax and Semantics of the input and output data and the Function of the AIW and the AIMS, and the Connections between and among the AIMS of an AIW.

In particular, an AIM is defined by its Function and data, but not by its internal architecture, which may be based on AI or data processing, and implemented in software, hardware or hybrid software and hardware technologies.

MPAI Standards are designed to enable a User to obtain, via standard protocols, an Implementation of an AIW and of the set of corresponding AIMS, and execute it in an AIF Implementation. The MPAI Store in *Figure 1* is an entity from which Implementations are downloaded. MPAI Standards assume that the AIF, AIW, and AIM Implementations may have been developed by independent implementers. A necessary condition for this to be possible, is that any AIF, AIW, and AIM implementations be uniquely identified. MPAI has appointed an ImplementerID Registration Authority (IIDRA) to assign unique ImplementerIDs (IID) to Implementers.<sup>1</sup>

A necessary condition to make possible the operations described in the paragraph above is the existence of an ecosystem composed of Conformance Testers, Performance Assessors, an instance of the IIDRA and of the MPAI Store. Reference [4] provides an informative example of such ecosystem.

<sup>1</sup> At the time of publication of this standard, the MPAI Store was assigned as the IIDRA.

The chapters and the annexes of this Technical Specification are Normative, unless they are labelled as Informative.

## 2 Scope of Standard

MPAI-NNW specifies methodologies to evaluate the following aspects of a neural network watermarking technology:

- The impact on the performance of a watermarked neural network and its inference.
- The ability of a neural network watermarking detector/decoder to detect/decode a payload when the watermarked neural network has been modified.
- The computational cost of injecting, detecting or decoding a payload in the watermarked neural network.

The standard assumes that:

- The neural network watermarking technology to be evaluated according to this standard is publicly available.
- The watermarking key is unknown during evaluation.
- The performance of the neural network watermarking technology does not depend on a specific key.

This Technical Specification has been developed by the MPAI Neural Network Watermarking Development Committee (NNW-DC). As the neural network watermarking area is fast-evolving, MPAI expects it will produce future MPAI-NNW versions providing methods to cope with technology evolution.

## 3 Terms and Definitions

The terms used in this standard whose first letter is capital have the meaning defined in *Table 1*.

*Table 1 – Table of terms and definitions*

<b>Term</b>	<b>Definition</b>
Computational cost	The cost of injecting, detecting or decoding a watermark in a neural network or its inference.
Imperceptibility	A difference in the performance of a neural network before and after the watermark embedding process.
Means	Procedure, tools, dataset or dataset characteristics used to evaluate one or more of Computational cost, Imperceptibility, or Robustness of a neural network watermarking technology.
Modification	The result of a simulated attack performed during Neural Network Watermarking testing.
Neural Network	or Artificial Neural Network, a set of interconnected information processing nodes whose connections are affected by Weights.
Neural Network Watermarking	The process of injecting a data payload in the Weights or the activation function of a Neural Network.
Parameter	A set of values characterizing the strength of a Modification.
Payload	The amount of information carried by the watermark.
Robustness	The ability of a watermarked neural network to withstand the impact of modifications in terms of detection and decoding capability.
Tester	The user who evaluates a neural network watermarking technology according to this Technical Specification.
Weight	The value by which the connection between two nodes of a Neural

## 4 Use cases (Informative)

This chapter provides an overview of possible use cases of MPAI-NNW together with the types of actors playing roles in them. These are provided for information and are not intended to restrict the scope of application of the standard.

The following use cases can relate to both watermarking the NN model or the NN inference:

- *Identify an NN*  
In this use case, the retrieved Payload conveys information about the NN itself.
- *Identify the actors of an NN*  
Actors are any of NN customer, NN end-user, NN owner, and NN watermarking provider.  
In this use case, the retrieved Payload conveys information about some or all of the following
- *Verify the integrity of an NN*  
In this use case, the Payload conveys information about the NN Model's integrity.
- *Assess the computational cost of injecting, detecting, and decoding a payload*

## 5 References

### 5.1 Normative references

MPAI-AIF normatively references the following documents:

1. MPAI; The MPAI Statutes; <https://mpai.community/statutes/>
2. MPAI; The MPAI Patent Policy; <https://mpai.community/about/the-mpai-patent-policy/>.
3. MPAI; Framework Licence of the Artificial Intelligence Framework Technical Specification (MPAI-AIF); <https://mpai.community/standards/mpai-aif/framework-licence/>

### 5.2 Informative references

4. Technical Specification: The Governance of the MPAI Ecosystem V1, 2021; <https://mpai.community/standards/mpai-gme/>

## 6 Imperceptibility evaluation

This chapter will deal with two cases:

- NNs for which the watermark is added after the NNs model was created.
- NNs for which the watermark is added during the training of the NNs model.

### 6.1 Watermark embedding is done after training

The Imperceptibility evaluation specifies the Means that enable a Tester to evaluate the differences in performance of a neural network before and after the watermark embedding process. There are two cases:

1. The NN has the input and output data format with specified semantics.
2. The input and output data format of the NN do not have specified semantics.

#### 6.1.1 Evaluation of an NN whose I/O data format has specified semantics

In this section, two actors are involved: the NN Watermarking provider requesting a Tester to evaluate the Imperceptibility performance of their watermarking technology.

The Tester shall adopt the following procedure:

1. Define a pair of training and testing datasets with a size with at least an order of magnitude more entries than trainable parameters.
2. Select:

- a. A set of  $M$  unwatermarked NNs trained on the training dataset.
- b.  $D$  data payloads corresponding to the pre-established payload size.
3. Apply the watermarking technology to the  $M$  NNs.
4. Process the training dataset and the  $D$  data payloads (if needed).
5. Feed the  $M$  unwatermarked NN with the test dataset
6. Measure the task-dependent quality of the produced inference.
7. Feed the  $M \times D$  watermarked NN with the same test dataset
8. Measure the task-dependent quality of the produced inference.
9. Provide the task-dependent quality of the produced inference measured in 6 and 7.

### 6.1.2 Evaluation of an NN whose I/O data format has no specified semantics

In this section, two actors are involved: the NN Watermarking provider requesting a Tester to evaluate the Imperceptibility performance of their watermarking technology.

The workflow of the process shall be the following:

1. Tester connects the NN to other NN until the input and output of the resulting configuration have input / output formats with specified semantics.
2. Tester applies all the steps in 6.1.1.

## 6.2 Watermark embedding is done during training

The Imperceptibility evaluation specifies the Means for evaluating the performance of a watermarked neural network. The workflow of the process shall evaluate the watermarked NN as an NN.

## 7 Robustness evaluation

The Robustness evaluation specifies the Means to enable a Tester to evaluate the robustness of the watermark against a set of modifications requested by one of the Actors.

The Tester evaluates the decoder and detector capability of a watermarking technology as specified in the following workflow:

1. Select:
  - a. A set of  $M$  unwatermarked NNs trained on the training dataset.
  - b.  $D$  data payloads corresponding to the pre-established payload size.
2. Apply the watermarking technology to the  $M$  NNs with the  $D$  data payloads
3. Produce a set of  $M \times (D + 1)$  modified NNs ( $M$  unwatermarked NNs and  $M \times D$  watermarked NNs), by applying one of the Modifications in Table 3 to a given Parameter value.
4. Evaluate the Robustness of the detector:
  - a. Apply the Watermark detector to any of the  $M \times (D + 1)$  NNs
  - b. Record the corresponding binary detection results (Yes – the mark is detected or No – the mark is not detected) – see Figure 7.
  - c. Label the Yes/No outputs of the Watermark detector as *true positive*, *true negative*, *false positive (false alarm)* and *false negative (missed detection)* according to the actual result – see Table 1.
  - d. Count the total number of false positives and the total number of false negatives.
5. Evaluate the Robustness of the decoder:
  - a. Apply the Watermark decoder to any of the  $M \times (D + 1)$  NNs
  - b. Compute a Distance between the outputs of the decoder and their corresponding original data payloads.
  - c. Compute the Symbol Error Rate (SER) for any of the  $M \times (D + 1)$  NNs, as the ratio of the distance to the size of the corresponding data payload.
  - d. Compute the average SER, as the average over the  $M \times (D + 1)$  SER values computed in the previous step.

6. Provide the average values over the total number of tests:
  - a. The ratio of the number of false positives to  $M \times (D + 1)$ ,
  - b. The ratio of the number of false negatives to  $M \times (D + 1)$ .
  - c. The  $M \times D$  number for tested NNs, and the average SER.
7. Repeat steps 3, 4, 5 and 6 for the requested number of Parameters values chosen in the ranges provided by Table 2.
8. Repeat steps 3, 4, 5, 6 and 7 for the requested set of Modifications chosen in the ranges provided by Table 2.

Table 2. List of modification with their parameters

Modification name	Parameter type	Parameter range
Modification	Parameter type	Parameter range
<b>Gaussian noise addition:</b> adding a zero-mean, $S$ standard deviation Gaussian noise to a layer in the NN model. This noise addition can be simultaneously applied to a subset of layers.	<ul style="list-style-type: none"> <li>- the layers to be modified by Gaussian noise</li> <li>- the ratio of <math>S</math> to standard deviation of the weights in the corresponding layer</li> </ul>	<ul style="list-style-type: none"> <li>- 1 to total number of layers</li> <li>- 0.1 to 0.3</li> </ul>
<b>L1 Pruning:</b> delete the $P\%$ of the smallest weights, irrespective of their layers.	<ul style="list-style-type: none"> <li>- the <math>P</math> percentage of the deleted weights</li> </ul>	<ul style="list-style-type: none"> <li>- 1% to 90%</li> <li>- 1% to 99.99% when aiming one layer</li> </ul>
<b>Random pruning:</b> delete $R\%$ of randomly selected weights, irrespective of their layers.	<ul style="list-style-type: none"> <li>- the <math>R</math> percentage of the deleted weights</li> </ul>	<ul style="list-style-type: none"> <li>- 1% to 10%</li> </ul>
<b>Quantizing:</b> reduce to $B$ the number of bits used to represent the weights by <ol style="list-style-type: none"> <li>1. reducing the number of bits based on a sequence of three operations: affine mapping from the weights interval to the <math>(0; 2^B - 1)</math></li> <li>2. rounding to the closest integer</li> <li>3. backward affine mapping towards the initial weights interval</li> </ol>	<ul style="list-style-type: none"> <li>- the layers to be modified by quantization</li> <li>- the value of <math>B</math></li> </ul>	<ul style="list-style-type: none"> <li>- 1 to total number of layers</li> <li>- 32 to 2</li> </ul>
<b>Fine tuning / transfer learning:</b> resume the training of the $M$ watermarked NNs submitted to test, for $E$ additional epochs.	<ul style="list-style-type: none"> <li>- ratio of <math>E</math> to the number of epochs in the initial training</li> </ul>	<ul style="list-style-type: none"> <li>- up to 0.5 time the total number of epochs</li> </ul>
<b>Knowledge distillation:</b> train a surrogate network using the	<ul style="list-style-type: none"> <li>- The structure of the architecture</li> </ul>	<ul style="list-style-type: none"> <li>- structures N</li> </ul>



inferences of the NN under test as training dataset	- The size of the dataset $D$ - The number of epochs $E$	- 10,000 to 1,000,000 - 1 to 100
<b>Watermark overwriting:</b> successively insert $R$ additional watermarks, with random payloads of the same size as the initial watermark	- $R$ number of watermarks successively inserted	- 2 to 4

## 8 Computational cost evaluation

The Computational cost evaluation specifies the Means that enable a Tester to evaluate the computational cost of:

- Injecting, in terms of memory footprint, time to process an epoch, and number of epochs necessary to insert the watermark.
- Detecting or decoding, in terms of memory footprint and time for the detector or the decoder to produce the expected result.

### 8.1 Computational cost of injecting a watermark

The Computational cost evaluation specifies the Means that enable a Tester to evaluate the computational cost of the injection using neural network watermarking method under testing.

The following four elements shall be used to characterize the injection process:

1. The memory footprint.
2. The time to execute the operation required by one epoch normalized according to the number of batches processed in one epoch.
3. In case of the injection is done concurrently with the training of the network, the number of epochs required to insert the watermark.
4. The time for the watermarked neural network to compute an inference.

The Tester shall evaluate the Computational cost of the injection according to the following workflow:

1. Define a pair of training and testing datasets with a size with at least an order of magnitude more entries than trainable parameters.
2. Select:
  - a. The training dataset (if needed).
  - b. A set of  $M$  unwatermarked NNs trained on the training dataset.
  - c.  $D$  data payloads corresponding to the pre-established payload size.
3. Apply the watermarking technology to the  $M$  NNs using the  $D$  data payloads.
4. Record the corresponding  $M \times D$  set of values characterizing the processing.
5. Provide the statistical average of the values over the total number of tests (*i.e.*  $M \times D$ ) for one of the informative Testing Environments of Table 3.

### 8.2 Computational cost of detecting/decoding

The MPAI Computational cost evaluation specifies the Means that enable a Tester to evaluate the computational cost of the detecting/decoding of a neural network watermarking methods.

We use the total duration and the memory footprint to characterize the detecting/decoding process. The Tester shall evaluate the Computational cost of detecting/decoding according to the following workflow:

1. Select a set of  $M$  unwatermarked NNs,  $D$  data payloads corresponding to the pre-established payload size and, if needed, the train dataset.
2. Apply the watermarking technology to the  $M$  NNs with the  $D$  data payloads
3. Evaluate the Robustness of the detector:
  - a. Apply the Watermark detector to any the  $M \times D$  NNs.
  - b. Record the corresponding  $M \times D$  set of values characterizing the processing.
4. Evaluate the Robustness of the decoder:
  - a. Apply the Watermark decoder to any the  $M \times D$  NNs.
  - b. Record the corresponding  $M \times D$  set of values characterizing the processing.
5. Provide the statistical average of the values over the total number of tests (over  $M \times D$ ) for one of the informative Testing Environments of Table 3.

*Table 3. Testing Environments (informative)*

<b>Type</b>	<b>Testing environment</b>
<b>S</b>	- Single GPU (16GB/6144 CUDA cores) - 8 cores CPU (2.6GHz)
<b>L</b>	- Double GPU (32GB/12288 CUDA cores) - 16 cores CPU (3.4GHz)

## Annex 1 MPAI-wide terms and definitions

The Terms used in this standard whose first letter is capital and are not already included in *Table 1* are defined in *Table 4*.

*Table 4 – MPAI-wide Terms*

<b>Term</b>	<b>Definition</b>
Access	Static or slowly changing data that are required by an application such as domain knowledge data, data models, etc.
AI Framework (AIF)	The environment where AIWs are executed.
AI Module (AIM)	A data processing element receiving AIM-specific Inputs and producing AIM-specific Outputs according to according to its Function. An AIM may be an aggregation of AIMs.
AI Workflow (AIW)	A structured aggregation of AIMs implementing a Use Case receiving AIW-specific inputs and producing AIW-specific outputs according to the AIW Function.
Application Standard	An MPAI Standard designed to enable a particular application domain.
Channel	A connection between an output port of an AIM and an input port of an AIM. The term “connection” is also used as synonymous.
Communication	The infrastructure that implements message passing between AIMs
Component	One of the 7 AIF elements: Access, Communication, Controller, Internal Storage, Global Storage, Store, and User Agent
Conformance	The attribute of an Implementation of being a correct technical Implementation of a Technical Specification.
Conformance Tester	An entity Testing the Conformance of an Implementation.
Conformance Testing	The normative document specifying the Means to Test the Conformance of an Implementation.
Conformance Testing Means	Procedures, tools, data sets and/or data set characteristics to Test the Conformance of an Implementation.
Connection	A channel connecting an output port of an AIM and an input port of an AIM.
Controller	A Component that manages and controls the AIMs in the AIF, so that they execute in the correct order and at the time when they are needed
Data Format	The standard digital representation of data.
Data Semantics	The meaning of data.
Ecosystem	The ensemble of actors making it possible for a User to execute an application composed of an AIF, one or more AIWs, each with one or more AIMs potentially sourced from independent implementers.
Explainability	The ability to trace the output of an Implementation back to the inputs that have produced it.
Fairness	The attribute of an Implementation whose extent of applicability can be assessed by making the training set and/or network open to testing for bias and unanticipated results.
Function	The operations effected by an AIW or an AIM on input data.
Global Storage	A Component to store data shared by AIMs.

Internal Storage	A Component to store data of the individual AIMS.
Identifier	A name that uniquely identifies an Implementation.
Implementation	<ol style="list-style-type: none"> <li>1. An embodiment of the MPAI-AIF Technical Specification, or</li> <li>2. An AIW or AIM of a particular Level (1-2-3) conforming with a Use Case of an MPAI Application Standard.</li> </ol>
Implementer	A legal entity implementing MPAI Technical Specifications.
ImplementerID (IID)	A unique name assigned by the ImplementerID Registration Authority to an Implementer.
ImplementerID Registration Authority (IIDRA)	The entity appointed by MPAI to assign ImplementerID's to Implementers.
Interoperability	The ability to functionally replace an AIM with another AIW having the same Interoperability Level
Interoperability Level	<p>The attribute of an AIW and its AIMS to be executable in an AIF Implementation and to:</p> <ol style="list-style-type: none"> <li>1. Be proprietary (Level 1)</li> <li>2. Pass the Conformance Testing (Level 2) of an Application Standard</li> <li>3. Pass the Performance Testing (Level 3) of an Application Standard.</li> </ol>
Knowledge Base	Structured and/or unstructured information made accessible to AIMS via MPAI-specified interfaces
Message	A sequence of Records transported by Communication through Channels.
Normativity	The set of attributes of a technology or a set of technologies specified by the applicable parts of an MPAI standard.
Performance	The attribute of an Implementation of being Reliable, Robust, Fair and Replicable.
Performance Assessment	The normative document specifying the Means to Assess the Grade of Performance of an Implementation.
Performance Assessment Means	Procedures, tools, data sets and/or data set characteristics to Assess the Performance of an Implementation.
Performance Assessor	An entity Assessing the Performance of an Implementation.
Profile	A particular subset of the technologies used in MPAI-AIF or an AIW of an Application Standard and, where applicable, the classes, other subsets, options and parameters relevant to that subset.
Record	A data structure with a specified structure
Reference Model	The AIMS and their Connections in an AIW.
Reference Software	A technically correct software implementation of a Technical Specification containing source code, or source and compiled code.
Reliability	The attribute of an Implementation that performs as specified by the Application Standard, profile and version the Implementation refers to, e.g., within the application scope, stated limitations, and for the period of time specified by the Implementer.
Replicability	The attribute of an Implementation whose Performance, as Assessed by a Performance Assessor, can be replicated, within an agreed level, by another Performance Assessor.
Robustness	The attribute of an Implementation that copes with data outside of the stated application scope with an estimated degree of confidence.

Scope	The domain of applicability of an MPAI Application Standard
Service Provider	An entrepreneur who offers an Implementation as a service (e.g., a recommendation service) to Users.
Standard	The ensemble of Technical Specification, Reference Software, Conformance Testing and Performance Assessment of an MPAI application Standard.
Technical Specification	(Framework) the normative specification of the AIF. (Application) the normative specification of the set of AIWs belonging to an application domain along with the AIMs required to Implement the AIWs that includes: <ol style="list-style-type: none"> <li>1. The formats of the Input/Output data of the AIWs implementing the AIWs.</li> <li>2. The Connections of the AIMs of the AIW.</li> <li>3. The formats of the Input/Output data of the AIMs belonging to the AIW.</li> </ol>
Testing Laboratory	A laboratory accredited to Assess the Grade of Performance of Implementations.
Time Base	The protocol specifying how Components can access timing information
Topology	The set of AIM Connections of an AIW.
Use Case	A particular instance of the Application domain target of an Application Standard.
User	A user of an Implementation.
User Agent	The Component interfacing the user with an AIF through the Controller.
Version	A revision or extension of a Standard or of one of its elements.
Zero Trust	A model of cybersecurity primarily focused on data and service protection that assumes no implicit trust.

## **Annex 2 Notices and Disclaimers Concerning MPAI Standards (Informative)**

The notices and legal disclaimers given below shall be borne in mind when [downloading](#) and using approved MPAI Standards.

In the following, “Standard” means the collection of four MPAI-approved and [published](#) documents: “Technical Specification”, “Reference Software” and “Conformance Testing” and, where applicable, “Performance Testing”.

### Life cycle of MPAI Standards

MPAI Standards are developed in accordance with the [MPAI Statutes](#). An MPAI Standard may only be developed when a Framework Licence has been adopted. MPAI Standards are developed by especially established MPAI Development Committees who operate on the basis of consensus, as specified in Annex 1 of the [MPAI Statutes](#). While the MPAI General Assembly and the Board of Directors administer the process of the said Annex 1, MPAI does not independently evaluate, test, or verify the accuracy of any of the information or the suitability of any of the technology choices made in its Standards.

MPAI Standards may be modified at any time by corrigenda or new editions. A new edition, however, may not necessarily replace an existing MPAI standard. Visit the [web page](#) to determine the status of any given published MPAI Standard.

Comments on MPAI Standards are welcome from any interested parties, whether MPAI members or not. Comments shall mandatorily include the name and the version of the MPAI Standard and, if applicable, the specific page or line the comment applies to. Comments should be sent to the [MPAI Secretariat](#). Comments will be reviewed by the appropriate committee for their technical relevance. However, MPAI does not provide interpretation, consulting information, or advice on MPAI Standards. Interested parties are invited to join MPAI so that they can attend the relevant Development Committees.

### Coverage and Applicability of MPAI Standards

MPAI makes no warranties or representations of any kind concerning its Standards, and expressly disclaims all warranties, expressed or implied, concerning any of its Standards, including but not limited to the warranties of merchantability, fitness for a particular purpose, non-infringement etc. MPAI Standards are supplied “AS IS”.

The existence of an MPAI Standard does not imply that there are no other ways to produce and distribute products and services in the scope of the Standard. Technical progress may render the technologies included in the MPAI Standard obsolete by the time the Standard is used, especially in a field as dynamic as AI. Therefore, those looking for standards in the Data Compression by Artificial Intelligence area should carefully assess the suitability of MPAI Standards for their needs.

IN NO EVENT SHALL MPAI BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO: THE NEED TO PROCURE SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR

TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE PUBLICATION, USE OF, OR RELIANCE UPON ANY STANDARD, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE AND REGARDLESS OF WHETHER SUCH DAMAGE WAS FORESEEABLE.

MPAI alerts users that practicing its Standards may infringe patents and other rights of third parties. Submitters of technologies to this standard have agreed to licence their Intellectual Property according to their respective Framework Licences.

Users of MPAI Standards should consider all applicable laws and regulations when using an MPAI Standard. The validity of Conformance Testing is strictly technical and refers to the correct implementation of the MPAI Standard. Moreover, positive Performance Assessment of an implementation applies exclusively in the context of the [MPAI Governance](#) and does not imply compliance with any regulatory requirements in the context of any jurisdiction. Therefore, it is the responsibility of the MPAI Standard implementer to observe or refer to the applicable regulatory requirements. By publishing an MPAI Standard, MPAI does not intend to promote actions that are not in compliance with applicable laws, and the Standard shall not be construed as doing so. In particular, users should evaluate MPAI Standards from the viewpoint of data privacy and data ownership in the context of their jurisdictions.

Implementers and users of MPAI Standards documents are responsible for determining and complying with all appropriate safety, security, environmental and health and all applicable laws and regulations.

#### Copyright

MPAI draft and approved standards, whether they are in the form of documents or as web pages or otherwise, are copyrighted by MPAI under Swiss and international copyright laws. MPAI Standards are made available and may be used for a wide variety of public and private uses, e.g., implementation, use and reference, in laws and regulations and standardisation. By making these documents available for these and other uses, however, MPAI does not waive any rights in copyright to its Standards. For inquiries regarding the copyright of MPAI standards, please contact the [MPAI Secretariat](#).

The Reference Software of an MPAI Standard is released with the [MPAI Modified Berkeley Software Distribution licence](#). However, implementers should be aware that the Reference Software of an MPAI Standard may reference some third-party software that may have a different licence.

## Annex 3 The Governance of the MPAI Ecosystem (Informative)

### Level 1 Interoperability

With reference to *Figure 1*, MPAI issues and maintains a Technical Specification – called MPAI-AIF – whose components are:

1. An environment called AI Framework (AIF) running AI Workflows (AIW) composed of inter-connected AI Modules (AIM) exposing standard interfaces.
2. A distribution system of AIW and AIM Implementation called MPAI Store from which an AIF Implementation can download AIWs and AIMs.

A Level 1 Implementation shall be an Implementation of the MPAI-AIF Technical Specification executing AIWs composed of AIMs able to call the MPAI-AIF APIs.

Implementers' benefits	Upload to the MPAI Store and have globally distributed Implementations of
	- AIFs conforming to MPAI-AIF.
	- AIWs and AIMs performing proprietary functions executable in AIF.
Users' benefits	Rely on Implementations that have been tested for security.
MPAI Store's role	- Tests the Conformance of Implementations to MPAI-AIF.
	- Verifies Implementations' security, e.g., absence of malware.
	- Indicates unambiguously that Implementations are Level 1.

### Level 2 Interoperability

In a Level 2 Implementation, the AIW shall be an Implementation of an MPAI Use Case and the AIMs shall conform with an MPAI Application Standard.

Implementers' benefits	Upload to the MPAI Store and have globally distributed Implementations of
	- AIFs conforming to MPAI-AIF.
	- AIWs and AIMs conforming to MPAI Application Standards.
Users' benefits	- Rely on Implementations of AIWs and AIMs whose Functions have been reviewed during standardisation.
	- Have a degree of Explainability of the AIW operation because the AIM Functions and the data Formats are known.
Market's benefits	- Open AIW and AIM markets foster competition leading to better products.
	- Competition of AIW and AIM Implementations fosters AI innovation.
MPAI Store's role	- Tests Conformance of Implementations with the relevant MPAI Standard.
	- Verifies Implementations' security.
	- Indicates unambiguously that Implementations are Level 2.

### Level 3 Interoperability

MPAI does not generally set standards on how and with what data an AIM should be trained. This is an important differentiator that promotes competition leading to better solutions. However, the performance of an AIM is typically higher if the data used for training are in greater quantity and more in tune with the scope. Training data that have large variety and cover the spectrum of all cases of interest in breadth and depth typically lead to Implementations of higher "quality".

For Level 3, MPAI normatively specifies the process, the tools and the data or the characteristics of the data to be used to Assess the Grade of Performance of an AIM or an AIW.

Implementers' benefits	May claim their Implementations have passed Performance Assessment.
------------------------	---



Users' benefits Get assurance that the Implementation being used performs correctly, e.g., it has been properly trained.

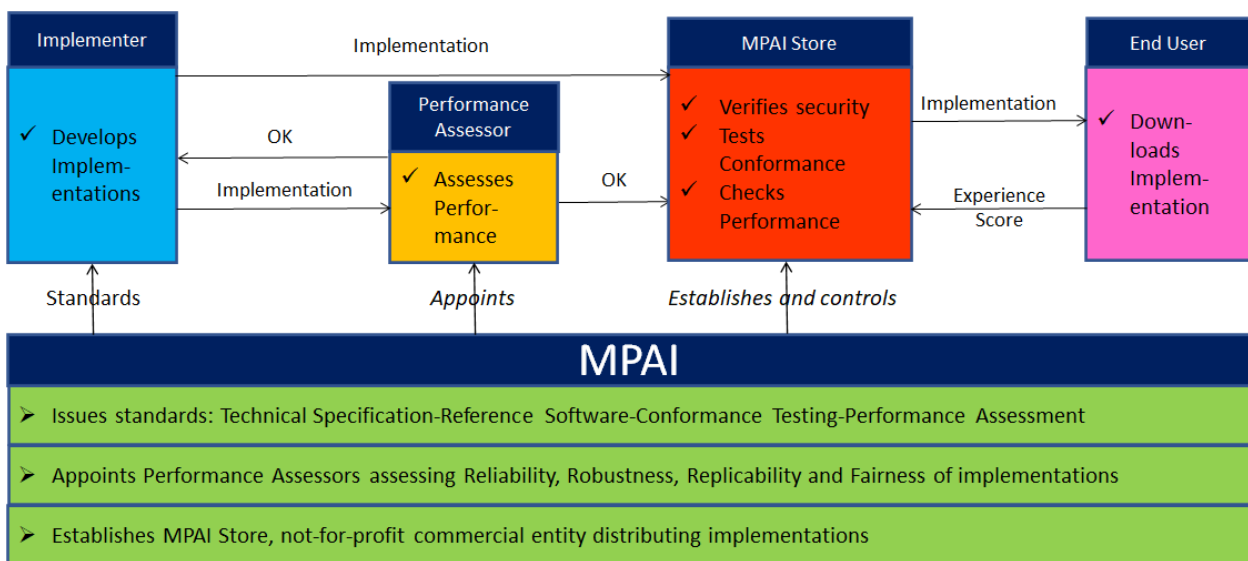
Market's benefits Implementations' Performance Grades stimulate the development of more Performing AIM and AIW Implementations.

MPAI Store's role - Verifies the Implementations' security  
- Indicates unambiguously that Implementations are Level 3.

### The MPAI ecosystem

The following *Figure 2* is a high-level description of the MPAI ecosystem operation applicable to fully conforming MPAI implementations as specified in the Governance of the MPAI Ecosystem Specification [4]:

1. MPAI establishes and controls the not-for-profit MPAI Store.
2. MPAI appoints Performance Assessors.
3. MPAI publishes Standards.
4. Implementers submit Implementations to Performance Assessors.
5. If the Implementation Performance is acceptable, Performance Assessors inform Implementers and MPAI Store.
6. Implementers submit Implementations to the MPAI Store
7. MPAI Store verifies security and Tests Conformance of Implementation.
8. Users download Implementations and report their experience to MPAI.



*Figure 2 – The MPAI ecosystem operation*

## **Annex 4 Patent declarations**

The MPAI Artificial Intelligence Framework (MPAI-AIF) Technical Specification has been developed according to the process outlined in the MPAI Statutes [1] and the MPAI Patent Policy [2].

The MPAI standardization process includes a step where by the secretariat issued a Call for Patent Declarations. The Table below will include information about the source of Patent Declarations that will be received in the future.

<b>Entity</b>	<b>Name</b>	<b>email address</b>

## Annex 5 Imperceptibility evaluation (informative)

### Classification task

The NN watermarking state of the art studies consider the classification as a predilection task. Within this task, the inference of a neural network belongs to a fix set of predefined classes.

To evaluate the impact of injecting a watermark in a classification NN:

- Probability of false alarm:  $P_{fa} = \frac{fp}{tp+fp+fn+tn}$  and Precision:  $\frac{tp}{tp+fp}$
- Probability of missed detection:  $P_{md} = \frac{fn}{tp+fp+fn+tn}$  and Recall:  $\frac{tp}{tp+fn}$

As these measures are based on binary classification problem, for multiclass classifiers the average for all classes shall be computed.

### Image/speech processing tasks

The inference of a neural network is a produced content. For example, a neural network for speech synthesis will return an artificial voice based on a text. Every qualitative/quantitative evaluation of a content can be use:

- Image: PSNR, SSIM, NCC, in addition to subjective test (e.g. as specified by ITU)
- Speech recognition: Word/Sentence error rate, Intent recognition rate, in addition to subjective test (e.g. as specified by ITU)

### Image semantic segmentation

The inference of a neural network is a semantic-labelled. To evaluate this method, we propose:

- Precision  $\frac{tp}{tp+fp}$
- Recall  $\frac{tp}{tp+fn}$
- Intersection over Union  $\frac{\text{Area of overlap}}{\text{Area of Union}}$