# OPEN SOURCE AI DEFINITION

*Online public townhall*

June 28, 2024

# Community agreements

- **One Mic, One Speaker** -- Please allow one person to speak at a time.
- **Take Space, Make Space** -- If you tend to talk more, we invite you to make space for others to share. If you tend not to share, we invite you to speak up.
- **Kindness** -- This work is hard, but we don't have to be. Gentleness and curiosity help. Those who use insults or hate speech will need to leave the meeting.
- **Forward Motion** -- We advance by focusing on what is possible in the moment and doing it. Obstacles are marked for later discussion, not used to stop the process. If we hit a boulder, we note it on the map and keep walking.  We'll come back and unearth it later on.
- **Solution-Seeking** -- This work is so complex that focusing on what won't work will stop it. Suggesting new ideas, options, and proposals is vulnerable, but crucial. All of us are needed to make this work.
- **Anything else?**

# OSI's objective for 2024
# **Open Source AI Definition**

Open Source AI Definition
**Current Version**
OSAID v.0.0.8

# Open Source AI Definition

v.0.0.8

## Preamble

## 4 Freedoms

## Legal Checklist

### Preamble

**Why we need Open Source Artificial Intelligence (AI)**

Open Source has demonstrated that massive benefits accrue to everyone when we remove the barriers to learning, using, sharing and improving software systems. These benefits are the result of using licenses that adhere to the Open Source Definition. The benefits can be summarized as autonomy, transparency, frictionless reuse, and collaborative improvement.

Everyone needs these benefits in AI. We need essential freedoms to enable users to build and deploy AI systems that are reliable and transparent.

### What is Open Source AI

- **Use** the system for any purpose and without having to ask for permission.
- **Study** how the system works and inspect its components.
- **Modify** the system for any purpose, including to change its output.
- **Share** the system for others to use with or without modifications, for any purpose.

Precondition to exercise these freedoms is to have access to the preferred form to make modifications to the system.

### Preferred form to make modifications to machine-learning systems

The preferred form of making modifications for a machine-learning Open Source AI must include:

- **Data information**: Sufficiently detailed information about the data used to train the system, so that a skilled person can recreate a substantially equivalent system using the same or similar data.
  - For example, if used, this would include the training methodologies and techniques, the training data sets used, information about the provenance of those data sets, their scope and characteristics, how the data was obtained and selected, the labeling procedures and data cleaning methodologies.
- **Code**: The source code used to train and run the system.
  - For example, if used, this would include code used for pre-processing data, code used for training, validation and testing, supporting libraries like tokenizers and hyperparameters search code, inference code, and model architecture.
- **Model**: The model parameters.
  - For example, this might include checkpoints from key intermediate stages of training as well as the final optimizer state.

### Checklist to evaluate machine learning systems

This checklist is based on the paper The Model Openness Framework: Promoting Completeness and Openness for Reproducibility, Transparency and Usability in AI published Mar 21, 2024.

**Table of default required components**

| Required components | Legal frameworks |
|---|---|
| **Data information** | |
| – Training methodologies and techniques | Available under OSD-compliant license |
| – Training data scope and characteristics | Available under OSD-compliant license |
| – Training data provenance (including how data was obtained and selected) | Available under OSD-compliant license |
| – Training data labeling procedures, if used | Available under OSD-compliant license |
| – Training data cleaning methodology | Available under OSD-compliant license |
| **Code** | |
| – Data pre-processing | Available under OSI-approved license |
| – Training, validation and testing | Available under OSI-approved license |
| – Inference | Available under OSI-approved license |
| – Supporting libraries and tools | Available under OSI-approved license |
| **Model** | |
| – Model architecture | Available under OSI-approved license |
| – Model parameters | Available under OSD-conformant terms |

The following components are not required as the preferred form of making modifications, but their inclusion in releases is appreciated.

| Optional components | Legal frameworks |
|---|---|
| **Data information** All data sets, including: | |
| – Training data sets | Available under OSD-compliant license |
| – Testing data sets | Available under OSD-compliant license |
| – Validation data sets | Available under OSD-compliant license |
| – Benchmarking data sets | Available under OSD-compliant license |
| – Data card | Available under OSD-compliant license |
| – Evaluation data | Available under OSD-compliant license |

Open Source AI Definition
**What We're Working On**
OSAID v.0.0.9

# Open Source AI Definition
## Preamble

v.0.0.9 plans

## Preamble

### Why we need Open Source Artificial Intelligence (AI)

Open Source has demonstrated that massive benefits accrue to everyone when you remove the barriers to learning, using, sharing and improving software systems. These benefits are the result of using licenses that adhere to the Open Source Definition. The benefits can be summarized as autonomy, transparency, frictionless reuse, and collaborative improvement.

Clarifying that the **recipients** of the freedoms are developers, deployers and end-users

# Open Source AI Definition

## Four Freedoms

v.0.0.9 plans

## What is Open Source AI

An Open Source AI is an AI system made available under terms that grant the freedoms to:

- **Use** the system for any purpose and without having to ask for permission.
- **Study** how the system works and inspect its components.
- **Modify** the system for any purpose, including to change its output.
- **Share** the system for others to use with or without modifications, for any purpose.

Clarifying that the four freedoms of open source AI are derived from the **Free Software Definition**

---

### Preamble

**Why we need Open Source Artificial Intelligence (AI)**

Open Source has demonstrated that massive benefits accrue to everyone when you remove the barriers to learning, using, sharing and improving software systems. These benefits are the result of using licenses that adhere to the Open Source Definition. The benefits can be summarized as autonomy, transparency, frictionless reuse, and collaborative improvement.

Everyone needs these benefits in AI. We need essential freedoms to enable users to build and deploy AI systems that are reliable and transparent.

**What is Open Source AI**

An Open Source AI is an AI system made available under terms that grant the freedoms to:

- **Use** the system for any purpose and without having to ask for permission.
- **Study** how the system works and inspect its components.
- **Modify** the system for any purpose, including to change its output.
- **Share** the system for others to use with or without modifications, for any purpose.

Precondition to exercise these freedoms is to have access to the preferred form to make modifications to the system.

**Preferred form to make modifications to machine-learning systems**

The preferred form of making modifications for a machine-learning Open Source AI must include:

- **Data information**: Sufficiently detailed information about the data used to train the system, so that a skilled person can recreate a substantially equivalent system using the same or similar data.
  - For example, if used, this would include the training methodologies and techniques, the training data sets used, information about the provenance of those data sets, their scope and characteristics, how the data was obtained and selected, the labeling procedures and data cleaning methodologies.
- **Code**: The source code used to train and run the system.
  - For example, if used, this would include code used for pre-processing data, code used for training, validation and testing, supporting libraries like tokenizers and hyperparameters search code, inference code, and model architecture.
- **Model**: The model parameters.
  - For example, this might include checkpoints from key intermediate stages of training as well as the final optimizer state.

**Checklist to evaluate machine learning systems**

This checklist is based on the paper The Model Openness Framework: Promoting Completeness and Openness for Reproducibility, Transparency and Usability in AI published Mar 21, 2024.

**Table of default required components**

| Required components | Legal frameworks |
| --- | --- |
| **Data information** | |
| – Training methodologies and techniques | Available under OSD-compliant license |
| – Training data scope and characteristics | Available under OSD-compliant license |
| – Training data provenance (including how data was obtained and selected) | Available under OSD-compliant license |
| – Training data labeling procedures, if used | Available under OSD-compliant license |
| – Training data cleaning methodology | Available under OSD-compliant license |
| **Code** | |
| – Data pre-processing | Available under OSI-approved license |
| – Training, validation and testing | Available under OSI-approved license |
| – Inference | Available under OSI-approved license |
| – Supporting libraries and tools | Available under OSI-approved license |
| **Model** | |
| – Model architecture | Available under OSI-approved license |
| – Model parameters | Available under OSD-conformant terms |

The following components are not required as the preferred form of making modifications, but their inclusion in releases is appreciated.

| Optional components | Legal frameworks |
| --- | --- |
| **Data information** All data sets, including: | |
| – Training data sets | Available under OSD-compliant license |
| – Testing data sets | Available under OSD-compliant license |
| – Validation data sets | Available under OSD-compliant license |
| – Benchmarking data sets | Available under OSD-compliant license |
| – Data card | Available under OSD-compliant license |
| – Evaluation data | Available under OSD-compliant license |

# Open Source AI Definition

## Four Freedoms

v.0.0.9 plans

## What is Open Source AI

An Open Source AI is an AI system made available under terms that grant the freedoms to:

- **Use** the system for any purpose and without having to ask for permission.
- **Study** how the system works and inspect its components.
- **Modify** the system for any purpose, including to change its output.
- **Share** the system for others to use with or without modifications, for any purpose.

---

### Preamble

#### Why we need Open Source Artificial Intelligence (AI)

Open Source has demonstrated that massive benefits accrue to everyone when you remove the barriers to learning, using, sharing and improving software systems. These benefits are the result of using licenses that adhere to the Open Source Definition. The benefits can be summarized as autonomy, transparency, frictionless reuse, and collaborative improvement.

Everyone needs these benefits in AI. We need essential freedoms to enable users to build and deploy AI systems that are reliable and transparent.

#### What is Open Source AI

An Open Source AI is an AI system made available under terms that grant the freedoms to:

- **Use** the system for any purpose and without having to ask for permission.
- **Study** how the system works and inspect its components.
- **Modify** the system for any purpose, including to change its output.
- **Share** the system for others to use with or without modifications, for any purpose.

Precondition to exercise these freedoms is to have access to the preferred form to make modifications to the system.

#### Preferred form to make modifications to machine-learning systems

The preferred form of making modifications for a machine-learning Open Source AI must include:

- **Data information**: Sufficiently detailed information about the data used to train the system, so that a skilled person can recreate a substantially equivalent system using the same or similar data.
  - For example, if used, this would include the training methodologies and techniques, the training data sets used, information about the provenance of those data sets, their scope and characteristics, how the data was obtained and selected, the labeling procedures and data cleaning methodologies.
- **Code**: The source code used to train and run the system.
  - For example, if used, this would include code used for pre-processing data, code used for training, validation and testing, supporting libraries like tokenizers and hyperparameters search code, inference code, and model architecture.
- **Model**: The model parameters.
  - For example, this might include checkpoints from key intermediate stages of training as well as the final optimizer state.
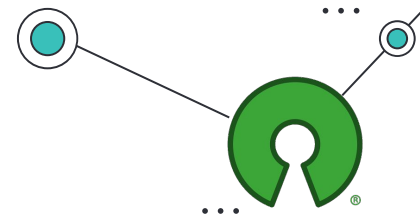
#### Checklist to evaluate machine learning systems

This checklist is based on the paper The Model Openness Framework: Promoting Completeness and Openness for Reproducibility, Transparency and Usability in AI published Mar 21, 2024.

**Table of default required components**

| Required components | Legal frameworks |
|---|---|
| **Data information** | |
| - Training methodologies and techniques | Available under OSD-compliant license |
| - Training data scope and characteristics | Available under OSD-compliant license |
| - Training data provenance (including how data was obtained and selected) | Available under OSD-compliant license |
| - Training data labeling procedures, if used | Available under OSD-compliant license |
| - Training data cleaning methodology | Available under OSD-compliant license |
| **Code** | |
| - Data pre-processing | Available under OSI-approved license |
| - Training, validation and testing | Available under OSI-approved license |
| - Inference | Available under OSI-approved license |
| - Supporting libraries and tools | Available under OSI-approved license |
| **Model** | |
| - Model architecture | Available under OSI-approved license |
| - Model parameters | Available under OSD-conformant terms |

The following components are not required as the preferred form of making modifications, but their inclusion in releases is appreciated.

| Optional components | Legal frameworks |
|---|---|
| **Data information** All data sets, including: | |
| - Training data sets | Available under OSD-compliant license |
| - Testing data sets | Available under OSD-compliant license |
| - Validation data sets | Available under OSD-compliant license |
| - Benchmarking data sets | Available under OSD-compliant license |
| - Data card | Available under OSD-compliant license |
| - Evaluation data | Available under OSD-compliant license |

---

Underlining that components and systems must be free from encumbrances that prevent any developer, deployer, or users from **exercising** those freedoms.

# Open Source AI Definition
## Preferred Form

v.0.0.9 plans



Adding definitions of…

… the "OSD **compliant**" requirement for data information...

…and the "OSD **conformant**" requirement for model parameters

..so legal requirements are clear for each component

---

**Preferred form to make modifications to machine-learning systems**

The preferred form of making modifications for a machine-learning Open Source AI must include:

- **Data Information**: Sufficiently detailed information about the data used to train the system, so that a skilled person can recreate a substantially equivalent system using the same or similar data.
  - For example, if used, this would include the training methodologies and techniques, the training data sets used, information about the provenance of those data sets, their scope and characteristics, how the data was obtained and selected, the labeling procedures and data cleaning methodologies.
- **Code**: The source code used to train and run the system.
  - For example, if used, this would include code used for pre-processing data, code used for training, validation and testing, supporting libraries like tokenizers and hyperparameters search code, inference code, and model architecture.
- **Model**: The model parameters.
  - For example, this might include checkpoints from key intermediate stages of training as well as the final optimizer state.

# Open Source AI Definition **Checklist**

v.0.0.9 plans



## Checklist to evaluate machine learning systems

This checklist is based on the paper The Model Openness Framework: Promoting Completeness and Openness for Reproducibility, Transparency and Usability in AI published Mar 21, 2024.

### Table of default required components

| Required components | Legal frameworks |
|---|---|
| **Data information** | |
| – Training methodologies and techniques | Available under OSD-compliant license |
| – Training data scope and characteristics | Available under OSD-compliant license |
| – Training data provenance (including how data was obtained and selected) | Available under OSD-compliant license |

Checklist will be a **separate** document and process and its components will be updated to follow the **Model Openness Framework** (MOF) precisely.

Open Source AI Definition
**System Validation**
OSAID v.0.0.8 (and soon v. 0.0.9)

# Validation Updates

Thanks to **Arctic** and **LLM360** for helping identify documentation!

| AI System | Meets OSAID requirements? | Notes |
|---|---|---|
| Name of system with link to its review sheet | Based on OSAID v. 0.0.8 and/or v.0.0.6 | Summary explanation of status (as of 6/11/24) |
| Arctic | Expect Yes | Verbal confirmation from Snowflake, which is adding legal documents to review sheet (6/3/24) |
| BLOOM | Confirmed No (license fails) | Usage restrictions in RAIL license |
| Falcon | Expect No | Documents on training methodologies and techniques and training, validation and testing are missing |
| Grok | Expect No | Very little public information on system |
| Llama 2 | Confirmed No | Data pre-processing + training, validation and testing code are not available |
| LLM360 | Expect Yes | Self-certified as compliant on the forum, awaiting addition of reviewable documents to their sheet |
| Mistral | Confirmed No | Some data information and code components missing, no training code available |
| OLMo | Expect Yes | Supporting libraries and tools unclear, but all other legal documentation is present |
| OpenCV | Unclear | Model requirement unclear because OpenCV does not store, but instead supports external deep learning frameworks |
| Phi-2 | Unclear | Data information, code, and model information missing |
| Poro | Unclear | Most review documentation not yet located; Located documentation meets OSAID requirements |
| Pythia | Confirmed Yes | Only non-alignment was absence of labeling documentation, which was not created. v 0.0.8 adds "if used" to requirement, resolving this |
| T5 | Expect Yes | Only possible restriction is in supporting libraries and tools because gcloud command requires special hardware. Hardware requirements are out of scope for the OSAID, so this is likely not a recognized restriction. |

# Open Source AI Definition
## What's Next?
## June - October 2024

- Complete validation phase
- Resolve comments, release v. 0.0.9 after validation
- Cut the release candidate with sufficient endorsement

# 2024 Timeline

| | System testing work stream |
|---|---|
| | Stakeholder consultation work stream |
| | Release schedule |

| February | June | July | August | September | October |
|---|---|---|---|---|---|
| **Call For Volunteers + Activity Feedback and Revision** | **Virtual System Review** | **Virtual System Review** | **Virtual System Review** | **Virtual System Review Ends** | |
| Bi-Weekly Virtual Public Townhalls | Bi-Weekly Virtual Public Townhalls | Townhalls + <br><br>**- OSPOs for Good** (NYC) <br>**- Sustain Africa** (virtual) | Townhalls + <br><br>**- AI-dev** (Hong Kong) | Townhalls + <br><br>**- Nerdearla** (Buenos Aires) | Townhalls + <br><br>**- All Things Open** (Raleigh) <br>**- Data Workshop** (Europe TBD) |
| Draft 0.0.5 | Draft 0.0.8 | Draft 0.0.9 | RC1 | RC1 | Stable Version |

# In-Person Meetings

| Region | Country | City | Conference | Date |
|---|---|---|---|---|
| North America | United States | Pittsburgh | ✔ **PyCon US** | May 17 |
| Europe | France | Paris | ✔ **OW2** | June 11 – 12 |
| North America | United States | New York | **OSPOs for Good** | July 9 – 11 |
| Africa | Virtual | Virtual | **Sustain Africa** | July 15 |
| Asia Pacific | China | Hong Kong | **AI_dev** | August 23 |
| Latin America | Argentina | Buenos Aires | **Nerdearla** | September 24 – 28 |
| Europe | TBD | TBD | **(data governance)** | October |
| North America | United States | Raleigh | **All Things Open** | Oct 27 – 29 |

# How to Participate :)



- Public **forum**: discuss.opensource.org

- Become an OSI **member**

  - Free or or full

  - SSO with other OSI websites

- Biweekly virtual **townhalls**… like this one!

- **Volunteer** to help with validation (email or DM Mer Joyce)

17

Q & A

# Thank you

We realize this is difficult work and we appreciate your help and openness in improving the definition.