

# OPEN SOURCE AI DEFINITION

*Online public townhall*

June 14, 2024

last updated: June 11, 2024 (MJ)

# ● 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

### 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, risk-reducing reuse, and collaborative improvement. Everyone reaps 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 the system

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/or on 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, **Model** 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 optimized 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 29, 2024.

##### Table of default required components

Required components	Legal frameworks
<b>Data information</b>	
- 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
<b>Code</b>	
- Data pre-processing	Available under GSI-approved license
- Training, validation and testing	Available under GSI-approved license
- Inference	Available under GSI-approved license
- Supporting libraries and tools	Available under GSI-approved license
<b>Model</b>	
- Model architecture	Available under OSD-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
<b>Data information</b> 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

## 4 Freedoms

## Legal Checklist

A teal circle is positioned on the left side of the slide, with a thin vertical line extending from the top to the bottom of the frame passing through its center.

Open Source AI Definition

## **Key Feedback**

OSAID v.0.0.8

# Open Source AI Definition Data Information

v.0.0.8

## Data Information

- Training methodologies and techniques
- Training data scope and characteristics
- Training data provenance (including how data was obtained and selected)
- Training data labeling procedures, if used
- Training data cleaning methodology

- Training data sets

## 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, risk-reducing reuse, and collaborative improvement. Everyone reaps 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.

Procedurally 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 on 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 tensors and hyperparameters search code, **Model** 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 29, 2024.

#### Table of default required components

Required components	Legal frameworks
<b>Data Information</b>	
- 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
<b>Code</b>	
- 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
<b>Model</b>	
- 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
- 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

Requiring only **data information**...

...instead of **training datasets** is the greatest point of debate now.

# Open Source AI Definition Other Components

v.0.0.8

## Code

- Data pre-processing

Available under OSI-approved license

- Data card

Available under OSD-compliant license

## 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, risk-reducing reuse, and collaborative improvement. Everyone reaps 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.

Provisional 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/or on 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 transformers and hyperparameters search code, **ModelCard** 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 29, 2024.

#### Table of default required components

Required components	Legal frameworks
<b>Data information</b>	
- Training methodologies and techniques	Available under OSI-compliant license
- Training data scope and characteristics	Available under OSI-compliant license
- Training data provenance (including how data was obtained and selected)	Available under OSI-compliant license
- Training data labeling procedures, if used	Available under OSI-compliant license
- Training data cleaning methodology	Available under OSI-compliant license
<b>Code</b>	
- 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
<b>Model</b>	
- 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
<b>Data information</b> AI data sets, including:	
- Training data sets	Available under OSI-compliant license
- Testing data sets	Available under OSI-compliant license
- Validation data sets	Available under OSI-compliant license
- Benchmarking data sets	Available under OSI-compliant license
- Data card	Available under OSI-compliant license
- Evaluation data	License

Others have proposed...

... removing **data pre-processing code** requirement if training data is not required.

...requiring a **model card**

... and **data card** to standardize system documentation.



# Open Source AI Definition Describing Legal Requirements

v.0.0.8

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

- Model parameters	Available under OSD-conformant terms
--------------------	--------------------------------------

## 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, risklessness, reuse, and collaborative improvement. Everyone reaps 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.

Provision 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/or 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 *scikit-learn* and *hyperparameters*, search code, *sklearn* 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

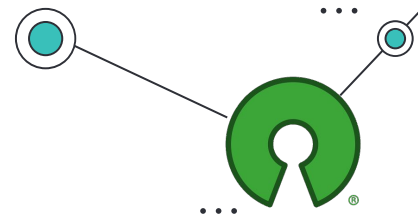
This checklist is based on the paper: The Model Openness Framework: Promoting Completeness and Openness for Reproducibility, Transparency and Usability in AI published Mar 21, 2024.

#### Table of default required components

Required components	Legal frameworks
<b>Data information</b>	
- 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
<b>Code</b>	
- 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
<b>Model</b>	
- 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
<b>Data information</b> AI 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



In contrast to the clear “OSI-approved” licenses available for code components...

... the “OSD compliant” requirement for data information...

...and “OSD conformant” requirement for model parameters have been challenging for reviewers to interpret.

- Preferred form to make modifications

- **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.

- Code**

The source code used to train and run the system.

- Model**

The model parameters  
*(weights and biases)*

- ## Data Information Explained

- The intention of *Data information* is to allow developers to **recreate** a substantially **equivalent system** using **the same or similar** data.
- Came out of the systems review process, with votes by volunteers.

## • Zooming in on the issues with datasets

- The Pile taken down after an alleged copyright infringement in the US. But legal in Japan. Maybe legal in EU
- DOLMA, initially had a restrictive license. Later switched to a permissive one. Suffers from the same legal uncertainties of the Pile, however the Allen Institute has not been sued, yet.
- Training techniques that preserve privacy like federated learning don't create datasets.

- Alternative proposals

- ● Use synthetic data: Experimental, unproven technology, limited to corner cases
- All their components must be “open source”: This integralism ignores that even the GNU project accepts system library exceptions and other compromises.



Open Source AI Definition

# **System Validation**

OSAID v.0.0.8

# Validation Reviewers

We're interested in reviewing about 10 AI systems self-described as open as part of this definition validation phase. Those marked (\*) have been reviewed in previous phases.

## 1. Arctic

1. **Jesús M. Gonzalez-Barahona** Universidad Rey Juan Carlos

## 2. BLOOM\*

2. **Danish Contractor** BLOOM Model Gov. Work Group
3. **Jaan Li** University of Tartu, One Fact Foundation

## 3. Falcon

1. **Casey Valk** Nutanix
2. **Jean-Pierre Lorre** LINAGORA, OpenLLM-France

## 4. Grok

1. **Victor Lu** independent database consultant
2. **Karsten Wade** Open Community Architects

## 5. Llama 2\*

1. **Davide Testuggine** Meta
2. **Jonathan Torres** Meta
3. **Stefano Zacchioli** Polytechnic Institute of Paris
4. **Victor Lu** independent database consultant

## 9. LLM360

5. **[Team member TBD]** LLM360

We will need an independent reviewer for LLM360

## 8. Mistral

1. **Mark Collier** OpenInfra Foundation
2. **Jean-Pierre Lorre** LINAGORA, OpenLLM-France
3. **Cailean Osborne** University of Oxford, Linux Foundation

## 7. OLMo

4. **Amanda Casari** Google
5. **Abdoulaye Diack** Google

## 8. OpenCV\*

1. **Rasim Sen** Oasis Software Technology Ltd.

## 9. Phi-2

6. **Seo-Young Isabelle Hwang** Samsung

## 10. Pythia\*

1. **Seo-Young Isabelle Hwang** Samsung
2. **Stella Biderman** EleutherAI
3. **Hailey Schoelkopf** EleutherAI
4. **Aviya Skowron** EleutherAI

## 11. T5

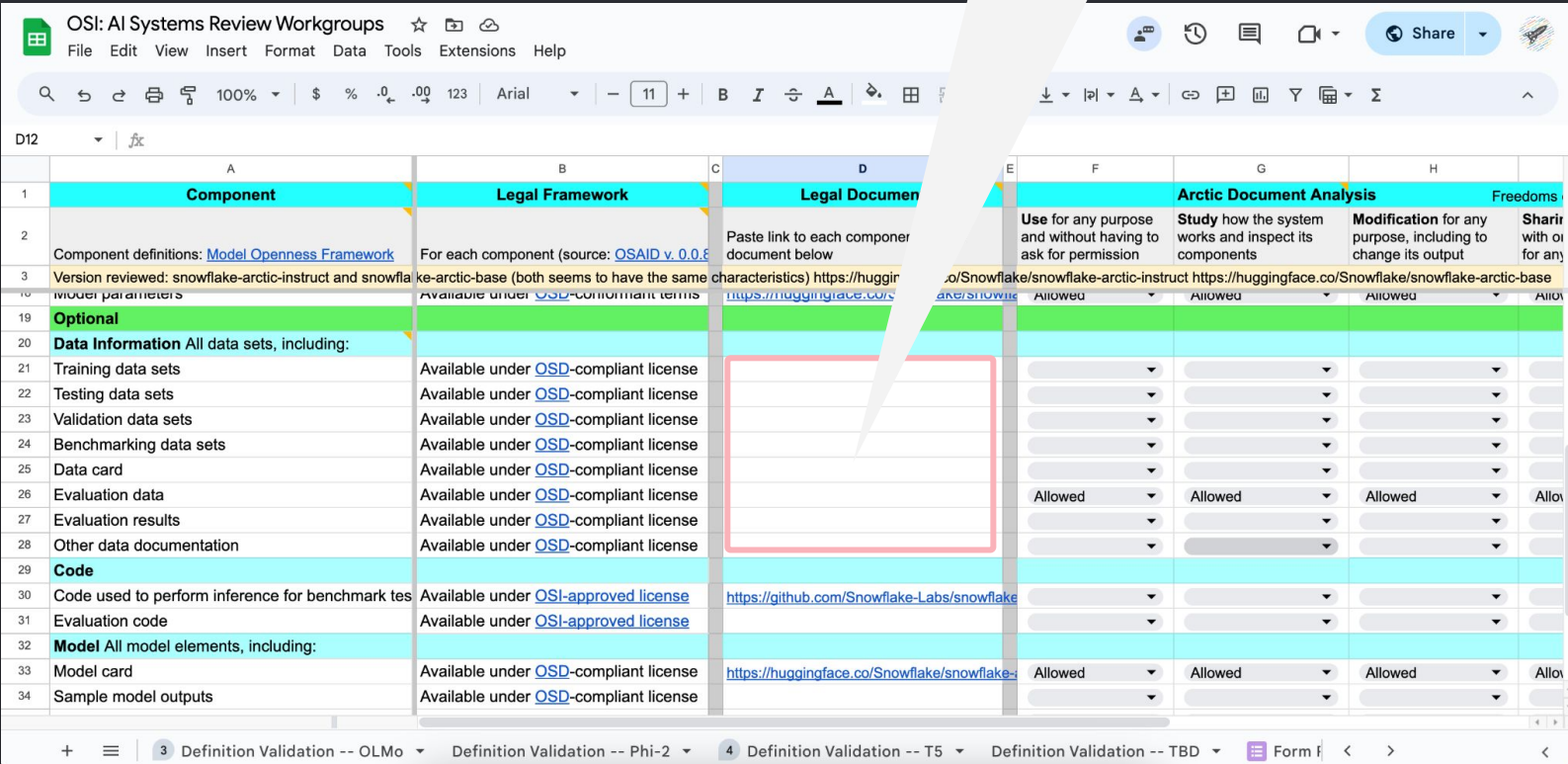
5. **Jaan Li** University of Tartu, One Fact Foundation

## Viking

6. **Merlijn Sebrechts** Ghent University

# Validation Challenges

It is hard for volunteer reviewers to find required documents independently..



OSI: AI Systems Review Workgroups

File Edit View Insert Format Data Tools Extensions Help

100% | \$ % .0 - .00 123 | Arial | 11 | B I | A |

D12

	A	B	C	D	E	F	G	H	I
1	<b>Component</b>	<b>Legal Framework</b>	<b>Legal Document</b>		<b>Arctic Document Analysis</b>				<b>Freedoms</b>
2	Component definitions: <a href="#">Model Openness Framework</a>	For each component (source: <a href="#">OSAIID v. 0.0.8</a> )	Paste link to each component document below		<b>Use for any purpose and without having to ask for permission</b>	<b>Study how the system works and inspect its components</b>	<b>Modification for any purpose, including to change its output</b>	<b>Share with others for any purpose</b>	
3	Version reviewed: snowflake-arctic-instruct and snowflake-arctic-base (both seems to have the same characteristics) <a href="https://huggingface.co/Snowflake/snowflake-arctic-instruct">https://huggingface.co/Snowflake/snowflake-arctic-instruct</a>	Available under <a href="#">OSD-compliant terms</a>	<a href="https://huggingface.co/Snowflake/snowflake-arctic-base">https://huggingface.co/Snowflake/snowflake-arctic-base</a>		Allowed	Allowed	Allowed	Allowed	Allowed
19	<b>Optional</b>								
20	<b>Data Information</b> All data sets, including:								
21	Training data sets	Available under <a href="#">OSD-compliant license</a>							
22	Testing data sets	Available under <a href="#">OSD-compliant license</a>							
23	Validation data sets	Available under <a href="#">OSD-compliant license</a>							
24	Benchmarking data sets	Available under <a href="#">OSD-compliant license</a>							
25	Data card	Available under <a href="#">OSD-compliant license</a>							
26	Evaluation data	Available under <a href="#">OSD-compliant license</a>			Allowed	Allowed	Allowed	Allowed	Allowed
27	Evaluation results	Available under <a href="#">OSD-compliant license</a>							
28	Other data documentation	Available under <a href="#">OSD-compliant license</a>							
29	<b>Code</b>								
30	Code used to perform inference for benchmark test	Available under <a href="#">OSI-approved license</a>	<a href="https://github.com/Snowflake-Labs/snowflake">https://github.com/Snowflake-Labs/snowflake</a>						
31	Evaluation code	Available under <a href="#">OSI-approved license</a>							
32	<b>Model</b> All model elements, including:								
33	Model card	Available under <a href="#">OSD-compliant license</a>	<a href="https://huggingface.co/Snowflake/snowflake-">https://huggingface.co/Snowflake/snowflake-</a>		Allowed	Allowed	Allowed	Allowed	Allowed
34	Sample model outputs	Available under <a href="#">OSD-compliant license</a>							

+ | 3 Definition Validation -- OLMO | Definition Validation -- Phi-2 | 4 Definition Validation -- T5 | Definition Validation -- TBD | Form | < >



# Validation Challenges

This meant a lot of the review analysis has been incomplete.

OSI: AI Systems Review Workgroups						
File Edit View Insert Format Data Tools Extensions Help						
100% 123 Arial 11 B I A						
D12						
1	Component	Legal Framework	Legal Document	Arctic Document	Arctic Document	Freedoms
2	Component definitions: <a href="#">Model Openness Framework</a>	For each component (source: <a href="#">OSAIID v. 0.0.8</a> )	Paste link to each component's legal document below	Use for any purpose and without having to ask for permission	Study how the system works and inspect its components	Modification for any purpose, including to change its output
3	Version reviewed: snowflake-arctic-instruct and snowflake-arctic-base (both seem to have the same characteristics)	<a href="#">https://huggingface.co/Snowflake/snowflake-arctic-instruct</a>	<a href="#">https://huggingface.co/Snowflake/snowflake-arctic-base</a>	Allowed	Allowed	Allowed
19	Optional	Available under <a href="#">OSD-compliant license</a>				
20	Data Information All data sets, including:					
21	Training data sets	Available under <a href="#">OSD-compliant license</a>				
22	Testing data sets	Available under <a href="#">OSD-compliant license</a>				
23	Validation data sets	Available under <a href="#">OSD-compliant license</a>				
24	Benchmarking data sets	Available under <a href="#">OSD-compliant license</a>				
25	Data card	Available under <a href="#">OSD-compliant license</a>				
26	Evaluation data	Available under <a href="#">OSD-compliant license</a>		Allowed	Allowed	Allowed
27	Evaluation results	Available under <a href="#">OSD-compliant license</a>				
28	Other data documentation	Available under <a href="#">OSD-compliant license</a>				
29	Code					
30	Code used to perform inference for benchmark tests	Available under <a href="#">OSI-approved license</a>	<a href="#">https://github.com/Snowflake-Labs/snowflake</a>			
31	Evaluation code	Available under <a href="#">OSI-approved license</a>				
32	Model All model elements, including:					
33	Model card	Available under <a href="#">OSD-compliant license</a>	<a href="#">https://huggingface.co/Snowflake/snowflake-</a>	Allowed	Allowed	Allowed
34	Sample model outputs	Available under <a href="#">OSD-compliant license</a>				

# Validation Solutions

Having the help of system creators to locate documents has been crucial.

Thank you, Arctic!

OSI: AI Systems Review Workgroups

File Edit View Insert Format Data Tools Extensions Help

100% 123 Arial 11 B I A

Component	Legal Framework	Legal Document	Arctic Document Analysis
Component definitions: <a href="#">Model Openness Framework</a>	For each component (source: <a href="#">OSAIID v. 0.0.8</a> )	Paste link to each component's legal document	Use for any purpose and Study how the system works Modification for any purpose
Version reviewed: snowflake-arctic-instruct and snowflake-arctic-instruct-2			
<b>Required</b>			
<b>Data Information</b>			
Training methodologies and techniques	Available under <a href="#">OSD-compliant license</a>	<a href="https://medium.com/snowflake/snowflake-arctic-instruct-2">https://medium.com/snowflake/snowflake-arctic-instruct-2</a>	Allowed
Training data scope and characteristics	Available under <a href="#">OSD-compliant license</a>	<a href="https://medium.com/snowflake/snowflake-arctic-instruct-2">https://medium.com/snowflake/snowflake-arctic-instruct-2</a>	Allowed
Training data provenance (including how data was collected)	Available under <a href="#">OSD-compliant license</a>	<a href="https://medium.com/snowflake/snowflake-arctic-instruct-2">https://medium.com/snowflake/snowflake-arctic-instruct-2</a>	Allowed
Training data labeling procedures, if used	Available under <a href="#">OSD-compliant license</a>	<a href="https://medium.com/snowflake/snowflake-arctic-instruct-2">https://medium.com/snowflake/snowflake-arctic-instruct-2</a>	Allowed
Training data cleaning methodology	Available under <a href="#">OSD-compliant license</a>	<a href="https://medium.com/snowflake/snowflake-arctic-instruct-2">https://medium.com/snowflake/snowflake-arctic-instruct-2</a>	Allowed
<b>Code</b>			
Data pre-processing	Available under <a href="#">OSI-approved license</a>	Something Snowflake is willing to share, but they haven't yet published this anywhere yet because no one has asked so far.	
Training, validation and testing	Available under <a href="#">OSI-approved license</a>	<a href="https://github.com/Snowflake-Labs/snowflake-arctic-instruct-2">https://github.com/Snowflake-Labs/snowflake-arctic-instruct-2</a>	Allowed
Inference	Available under <a href="#">OSI-approved license</a>	<a href="https://github.com/Snowflake-Labs/snowflake-arctic-instruct-2">https://github.com/Snowflake-Labs/snowflake-arctic-instruct-2</a>	Allowed
Supporting libraries and tools	Available under <a href="#">OSI-approved license</a>		Allowed
<b>Model</b>			
Model architecture	Available under <a href="#">OSI-approved license</a>	<a href="https://huggingface.co/Snowflake/snowflake-arctic-instruct-2">https://huggingface.co/Snowflake/snowflake-arctic-instruct-2</a>	Allowed
Model parameters	Available under <a href="#">OSD-conformant terms</a>	<a href="https://huggingface.co/Snowflake/snowflake-arctic-instruct-2">https://huggingface.co/Snowflake/snowflake-arctic-instruct-2</a>	Allowed
<b>Optional</b>			

12 Definition Validation -- Arctic 8 Definition Validation -- Falcon 4 Definition Validation -- Grok Definition Validation -- LLM360

# Validation Expectations

Given current system information, our expected review results are as follows. If we are missing information, please let us know.

AI System	Meets OSAID requirements?	Notes	Confirmed Yes	Expect Yes	Unclear	Expect No	Confirmed No
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)	Review complete + all required components present	Review incomplete + all required components expected to be present	Incomplete data, cannot ascertain openness	Review incomplete + some or all required components expected to be absent	Review complete + some/all required components absent
<a href="#">Arctic</a>	Expect Yes	Verbal confirmation from Snowflake, which is adding legal documents to review sheet (6/3/24)					
<a href="#">BLOOM</a>	Confirmed No	Usage restrictions in RAIL license					
<a href="#">Falcon</a>	Expect No	Documents on training methodologies and techniques and training, validation and testing are missing					
<a href="#">Grok</a>	Expect No	Very little public information on system					
<a href="#">Llama 2</a>	Confirmed No	Data pre-processing + training, validation and testing code are not available					
<a href="#">LLM360</a>	Expect Yes	Self-certified as compliant on the forum, awaiting addition of reviewable documents to their sheet					
<a href="#">Mistral</a>	Confirmed No	Some data information and code components missing, no training code available					
<a href="#">OLMo</a>	Expect Yes	Supporting libraries and tools unclear, but all other legal documentation is present					
<a href="#">OpenCV</a>	Unclear	Model requirement unclear because OpenCV does not store, but instead supports external deep learning frameworks					
<a href="#">Phi-2</a>	Unclear	Data information, code, and model information missing					
<a href="#">Pythia</a>	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					
<a href="#">T5</a>	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?**

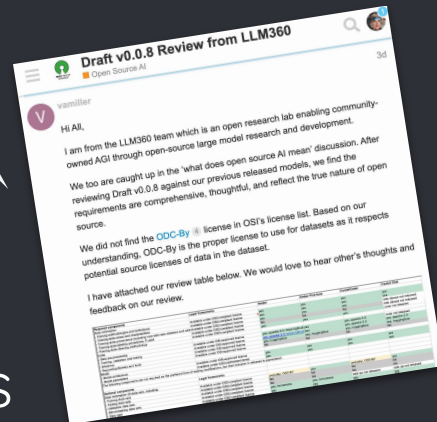
June - October 2024

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

# Complete the Validation Phase

1. Reach out to **AI system creators** to fill in the blanks on their own systems by pointing us to correct documentation
2. Invite **volunteers** to also help us fill in these blanks

Thanks,  
LLM360!

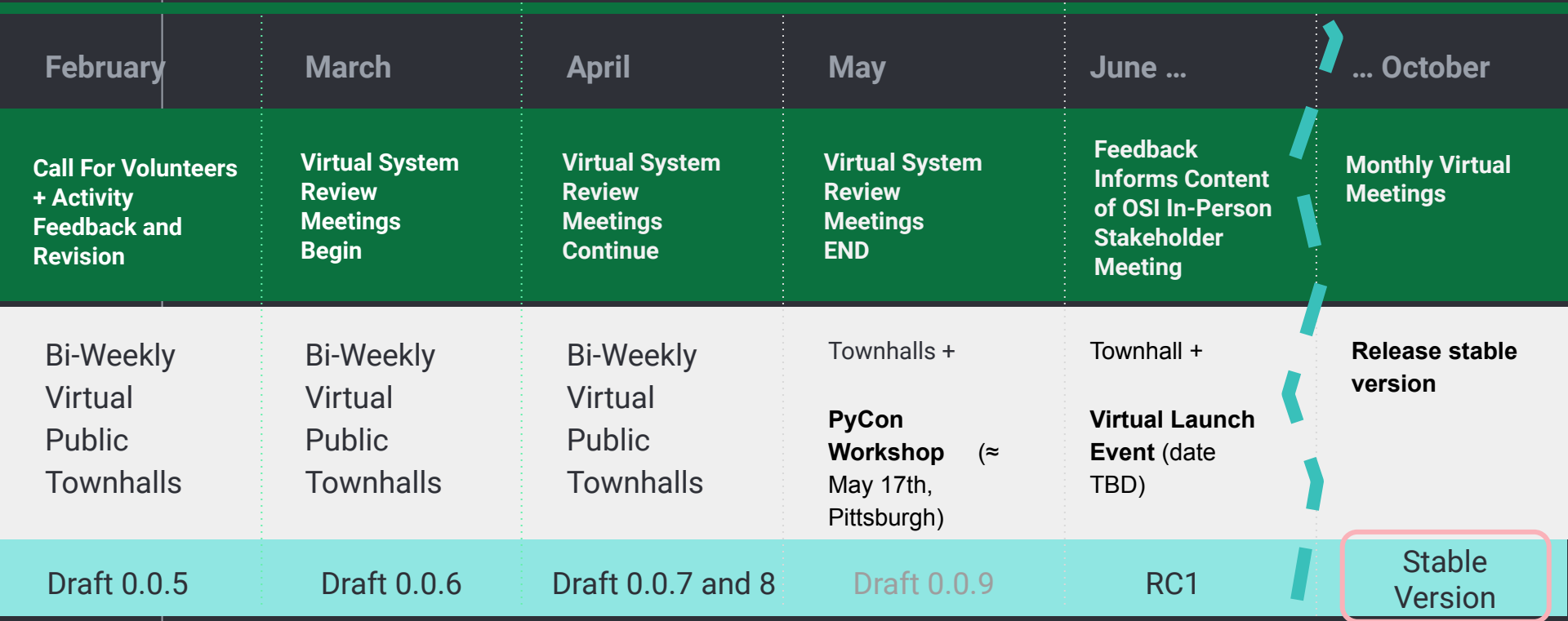


# 2024 Timeline

System testing work stream

Stakeholder consultation work stream

Release schedule

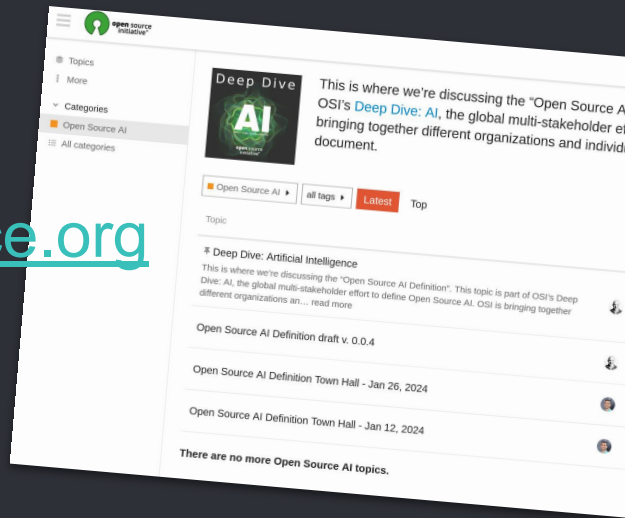


# In-Person Meetings

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	OSP0s for Good	July 9 - 11
Africa	Virtual	Virtual	Sustain Africa	July
Asia Pacific	China	Hong Kong	AI_dev	August 23
Latin America	Argentina	Buenos Aires	Nerdearla	September
Europe	TBD	TBD	(data governance)	October
North America	United States	Raleigh	All Things Open	Oct 27 - 29

# ● Participation Options

- Public forum: [discuss.opensource.org](https://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)







Q & A



## Thank you

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