# AI Alliance

## Open Trusted Data Initiative (OTDI)

https://the-ai-alliance.github.io/open-trusted-data-initiative/

June 20, 2025

thealliance.ai

# Focus Areas & Mission

Represents the investment priorities for the AI Alliance

*Member organizations have the choice to take part in one or more of these six focus areas and the agility to shift participation based on their interest and priorities.*

## Skills & Education

Support global AI skills building, education, and exploratory research.

## Trust & Safety

Create benchmarks, tools, and methodologies to ensure and evaluate high-quality and safe AI.

## Applications & Tools

Build and advance efficient and capable software frameworks for model builders and developers.

## HW Enablement

Foster a vibrant AI hardware accelerator ecosystem through SW.

## Foundation Models & Data

Enable an ecosystem of open foundation models and datasets for diverse modalities.

## Advocacy

Advocate for regulatory policies that create a healthy open ecosystem for AI.

AI Alliance

# Open Trusted Data Initiative (OTDI)

**Problem Statement**

Can I trust the datasets used for AI training, tuning, RAG, etc.?

○ Where did the data come from?

○ Is that OSS license (e.g., Apache) valid for *all* the data in the dataset?

○ What content is in the dataset? Copyrighted material, propaganda, ...?

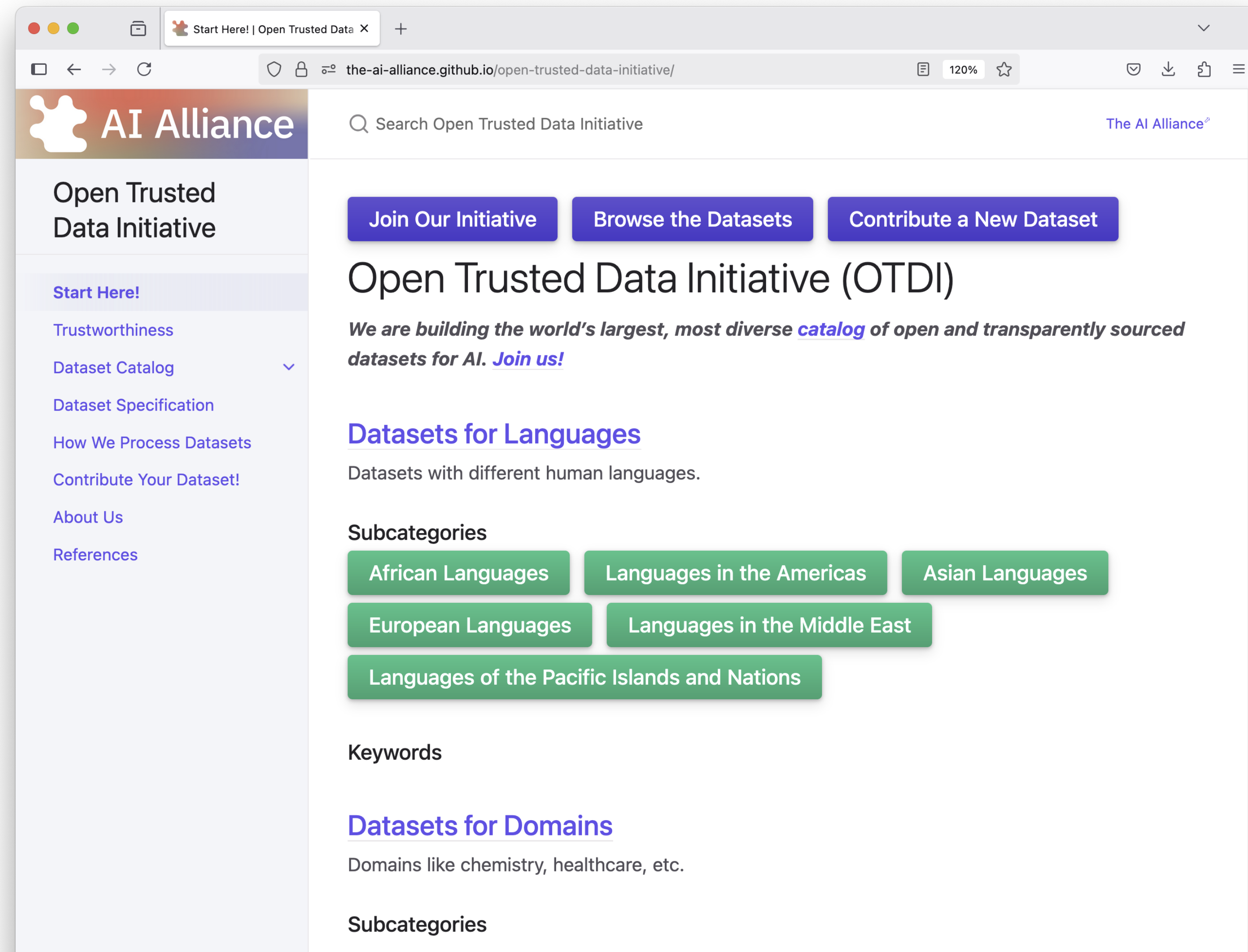○ How do I find datasets that meet my requirements?

An Initiative of <u>Focus Area 5: Foundation Models and Datasets</u>

# Open Trusted Data Initiative (OTDI)

## **What We Are Building**

We are building:

- **Definition of "open data":** Open-access, proper governance, clear provenance.

- **A dataset catalog:** Help users find suitable datasets.

- **Diversity of datasets:** Multilingual, multimodal, for domains: time series, science fields, industries, specific use cases.

- **Reusable tools:** Use our tools for your own needs.

An Initiative of <u>Focus Area 5: Foundation Models and Datasets</u>

<u>https://the-ai-alliance.github.io/open-trusted-data-initiative/</u>

# Open Trusted Data Initiative (OTDI)

## Definition of "open data"

What should *open data* mean?

- **Clear license:** *Every datum* is covered by an open-use license.

- **Clear provenance:** The origin of the dataset and its history are known.

- **Effective governance:** From creation, the dataset has been carefully managed.

Help us define what "open" means!

An Initiative of <u>Focus Area 5: Foundation Models and Datasets</u>

<u>https://the-ai-alliance.github.io/open-trusted-data-initiative/</u>

# Open Trusted Data Initiative (OTDI)

## Validation Tools

https://the-ai-alliance.github.io/open-trusted-data-initiative/

Verify metadata matches data:

- **Conformance checking:** Licenses are valid, contents match metadata, ...

- **Dataset transformations:** Create new datasets from other datasets.

Help us build these pipelines!

An Initiative of Focus Area 5: Foundation Models and Datasets

# Open Trusted Data Initiative (OTDI)

## A Dataset Catalog

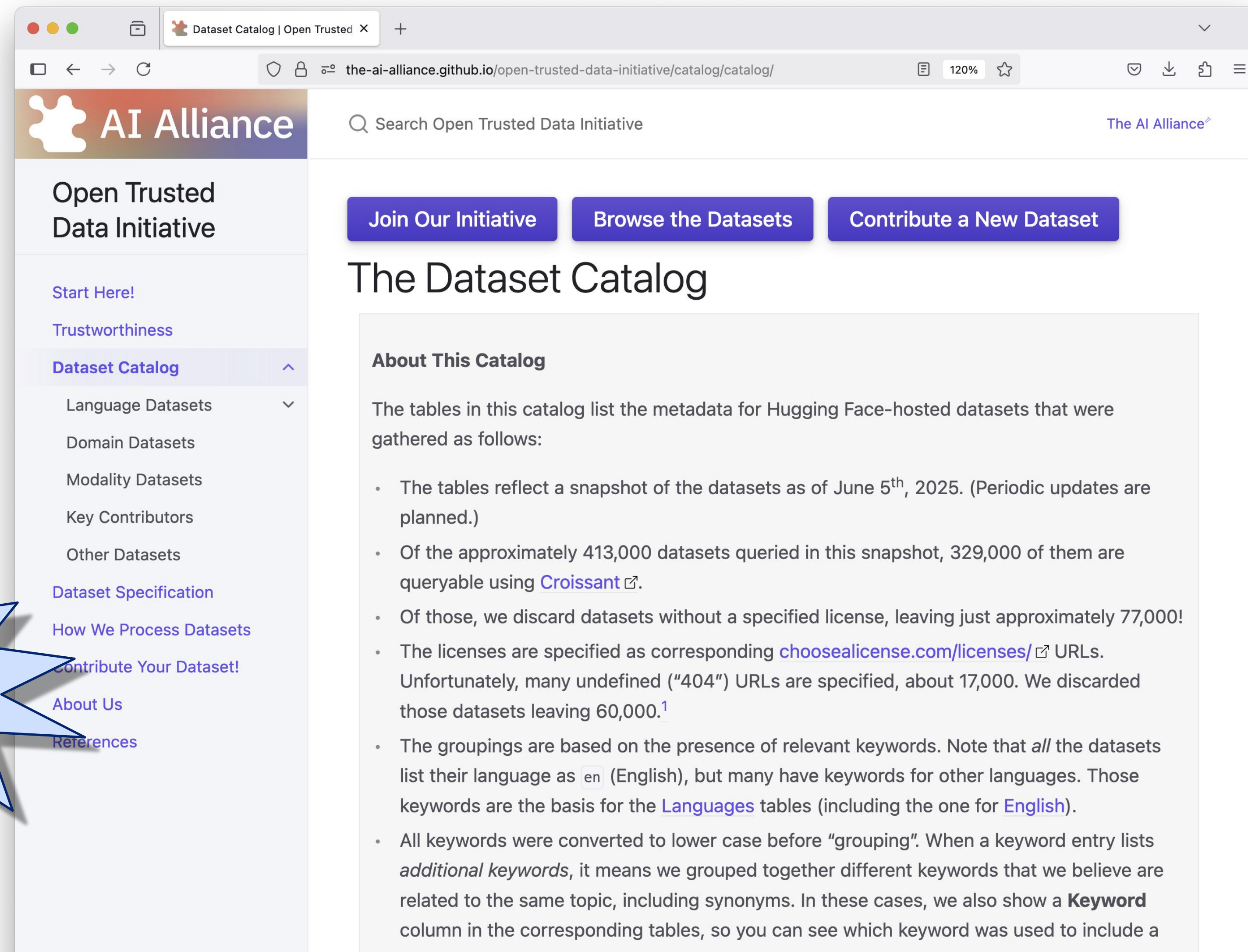https://the-ai-alliance.github.io/open-trusted-data-initiative/

Where are the open datasets?

- **Searchable:** By license, target application, domain, use case, ...

- **Track evolving datasets:** Updates *metadata* (Croissant) from known sources.

- **Many sources:** Hugging Face and others.

An Initiative of Focus Area 5: Foundation Models and Datasets



Help us build the catalog!

# Open Trusted Data Initiative (OTDI)

## Diversity of Datasets

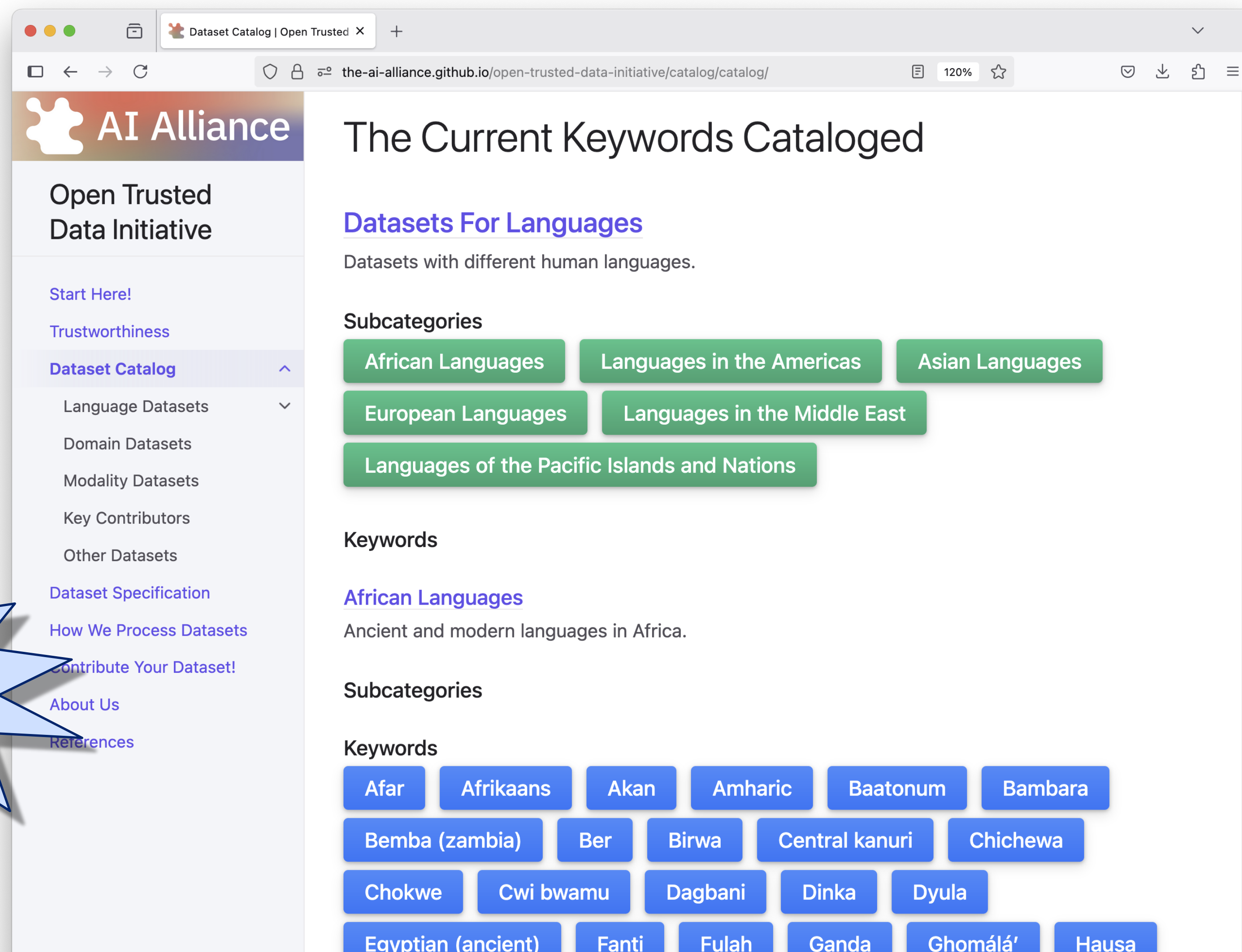https://the-ai-alliance.github.io/open-trusted-data-initiative/

What kinds of datasets?

- **For models:** LLM training, tuning, evaluation, ...

- **For domains:** Time-series, scientific discovery, vertical industries.

- **Multiple modalities:** Language, images, video, classification, ...

Use your domain expertise to create datasets

An Initiative of Focus Area 5: Foundation Models and Datasets

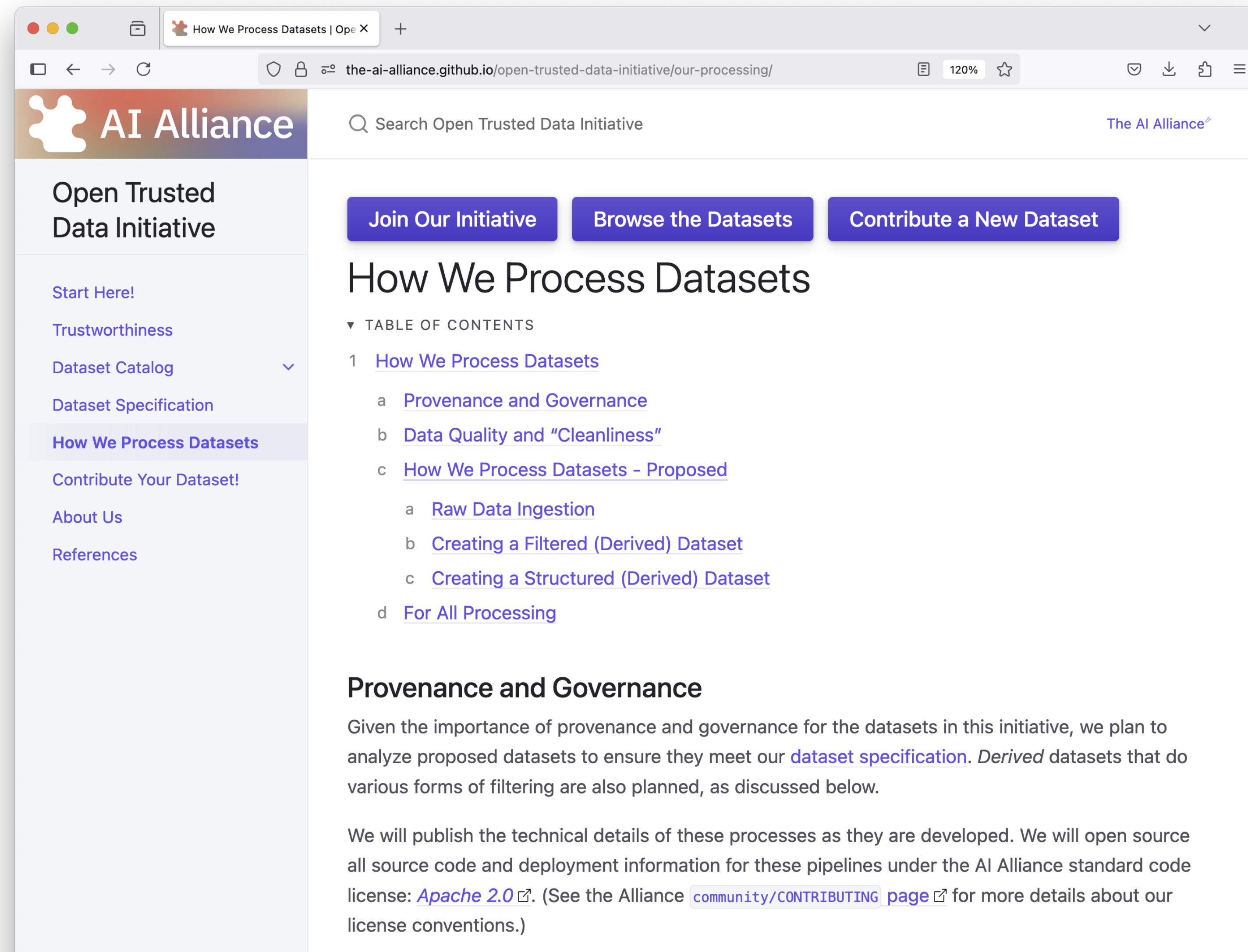# Open Trusted Data Initiative (OTDI)

## Reusable Tools

**https://the-ai-alliance.github.io/open-trusted-data-initiative/**

Use our tools for your own datasets.

○ **Conformance checking:** Licenses, contents, ...

○ **Dataset transformations:** Create new datasets from other datasets.

○ **Catalog:** Collect metadata for datasets and make it accessible.

An Initiative of Focus Area 5: Foundation Models and Datasets

# Open Trusted Data Initiative (OTDI)

**Please Join Us!!**

https://the-ai-alliance.github.io/open-trusted-data-initiative/

Help us help you...

○ Contribute your datasets

○ Help us build datasets for underserved languages and domains

○ Help us define and implement our standards and tools

Dean Wampler
dwampler@thealliance.ai
https://thealliance.ai