

Data Flywheels and Public AI

Nicholas Vincent

2025-08-04

Table of contents

Preface	4
1 Introduction	5
1.1 Definitions	5
1.2 Core Principles	6
2 The Serverless + Git MVP	7
2.1 Overview	7
2.2 Advantages of a Serverless + Git approach	8
2.3 Disadvantages (and what to watch)	9
2.4 Other reading:	10
3 Data Flow and Data Retention	11
4 Glossary of Defined Terms (for this Chapter)	12
4.1 1) System scope & components	13
4.2 2) What data is produced & when	14
4.2.1 2.1 Open WebUI (no account required, but optional and recommended)	14
4.2.2 2.2 Open WebUI Accounts (optional)	15
4.2.3 2.3 Model Gateway (to providers)	15
4.2.4 2.4 Opt-in Data Flywheel	15
4.3 3) Event timeline (how data flows)	16
4.4 4) Licensing & Preference Signals (Beta)	16
4.5 5) Retention schedule	17
4.6 6) Distribution & access control	18
4.6.1 6.1 Hugging Face (gated)	18
4.6.2 6.2 Flywheel Static Site (public)	18
4.7 7) Provider transparency (no guarantees)	19
4.8 8) Contributor rights & controls	19
4.9 9) Roadmap	19
4.10 10) Contacts	20
5 Rationale (and other flywheel variants)	21
5.1 Purpose of this section	21
5.2 More on all the other approaches we could've taken	21
5.2.1 Where does the "final" data live?	21

5.2.2	When is the user prompted to contribute?	21
5.2.3	What information object is created?	22
5.2.4	When is the data processed?	22
5.2.5	How is the data accessed?	22
5.2.6	How much friction is acceptable?	22
5.3	Some Categories of Architectural Models	23
5.3.1	Standard “PrivateCo” Web App	23
5.3.2	Git/Wiki Platform	23
5.3.3	Direct Telemetry	24
5.3.4	Hybrid Model	24
5.3.5	Serverless + Git Platform	24
5.3.6	Federated Learning Model	24
5.3.7	Browser Extension (just a data ingestion choice: can be used with various backend choices above)	25
5.3.8	P2P Network Model	25
5.4	Terse decision Matrix for the above dimensions	25
5.5	Scenario Walkthroughs: A Practical Comparison	26
5.5.1	Scenario A: User marks a chat as “Good” – when does processing happen?	27
5.5.2	Scenario B: User corrects a factual error	27
5.5.3	Scenario C: Accessing the contributed data	27
5.6	Frontier approaches data cooperatives, federated learning, and more	27
6	Appendix 1: LLM Data Schemas	29
7	Appendix 2 — Preference Signals for AI Data Use (CC signals + IETF AI Preferences)	32
8	Appendix 3: Example Legal Terms	34
8.1	Opt-in Data Flywheel — Legal Terms (Draft)	34
8.2	Frontend Instance	36

Preface

This is a “mini-book” that discusses “public AI flywheels”: software meant to enable people to opt-in to contribute data towards “public AI” causes. The goal of this book is to support efforts build a transparent, people-centric data collection ecosystem that supports the evaluation and training of public-benefit AI models. More frankly, this is way to organize some design notes, practical documentation that’s out of scope for a single example projec’s repo, and longer abstract writing on the topic.

This document is organized as such:

- In Section 1 (current section), we first we describe how this document is organized and introduce the Public AI Flywheel concept.
 - Section 1.2 describes “core principles” for a public AI data flywheel
- In Section 2, we discuss one particular implementation of a Minimum Viable Product (MVP) opt-in flywheel meant to accompany a “public AI interface” (hosted interface software that hits various endpoints for “public AI models”) that uses a “serverless” app + Git backend approach
 - This MVP focuses on collecting two high-signal data types: exports of “good chats” and “fail chats.” This data provides immediate value for model evaluation and, at scale, can be used for fine-tuning. Importantly, collecting a list of good and bad chats is also immediately fun, so contributors can get some value before we reach a threshold of data volume needed to construct a full benchmark or dataset. We expect key ideas discussed in this doc, and concretized in this project, to generalize to other data types.
- In Section 3, we provide details on the data retention policy for the Public AI Flywheel and the data policy for the full public AI interface pipeline: from model endpoints to OpenWebUI interface to flywheel platform.
- In Section 4, we explore the design space in much more detail, describing many other ways that flywheels could be implemented (with the hope of providing a helpful starting point, but also providing more rationale around our MVP)

1 Introduction

1.1 Definitions

What is a data flywheel? Nvidia gives us [this](#) definition: “A data flywheel is a feedback loop where data collected from interactions or processes is used to continuously refine AI models.” Others have also written on data flywheels (see e.g. a number of helpful blogs from [Roche and Sasson](#), [Liu](#) and [Del Balso](#)).

What is public AI? The public AI network gives us [this](#): AI with

“Public Access – Certain capabilities are so important for participation in public life that access to them should be universal. Public AI provides affordable access to these tools so that everyone can realize their potential.” “Public Accountability – Public AI earns trust by ensuring ultimate control of development rests with the public, giving everyone a chance to participate in shaping the future.” “Permanent Public Goods – Public AI is funded and operated in a way to maintain the public goods it produces permanently, enabling innovators to safely build on a firm foundation.”

For more on the public AI concept, see also Mozilla’s [work](#) in this space and several workshop papers and preprints (from [RegML 2023](#) at NeurIPS, [CodeML 2025](#) at ICML).

Our focus in this mini-book is in building public AI flywheels. In order to achieve public access and accountability, public AI systems must also face some unique challenges around the implementation of data flywheels – they may not be able to do what private orgs can do.

We are trying to create a feedback loop to improve AI: but we want to start from a position of accessibility (including providing an accessible explanation of exactly what happens to data) and accountability (so people have real agency over data).

Of course, it’s worth noting that a given public could deliberate and make a collective decision that they prefer an more “traditional approach” to data. Here, we are taking the stance that it’s best to start from a position of leaning heavily towards an opt-in approach (minimize usage and retention of data; data that is used in the flywheel to train AI should be provided via an opt-in by highly informed users).

1.2 Core Principles

Translating the core principles of public AI to the data flywheel domain (and to data strategy more generally), we aim to adhere to the following non-negotiable principles:

- **Transparency & Informed Consent:** Users must be fully informed about the model, its developer, and the ramifications of their contribution. A detailed FAQ and a clear consent module are required before any data is shared. To some extent, true transparency and informed consent requires to active expenditure of resources to improve the public’s AI literacy. We need systems that really do inform people (luckily, that’s something it seems like AI can be good for).
- **User Control & Data Rights** The system must empower users with control over their data, mirroring GDPR principles. This includes the right to access ([\\$Art. 15\\$](#)), rectify ([\\$Art. 16\\$](#)), erase ([\\$Art. 17\\$](#)), and port their data ([\\$Art. 20\\$](#)). Key exemplar: [Mozilla Common Voice](#).
 - We note that user control and data rights sometimes conflict with a “fully open” ethos; we will attempt to mitigate these tensions to the best extent possible!
- **Privacy and options for pseudonymity:** To the extent possible, we believe it is valuable to offer people to ability to contribute to data flywheel with their “real account” attached (to earn credit and reputation), but also to maintain an option to use “light-auth” system, requiring only an email (or similar) for authentication via passwordless magic links. If users take this approach, public contributions will be attributed to a randomly generated handle (e.g., “silver-badger-81”) to protect people’s identity. However, in our MVP (discussed in the next chapter) a HuggingFace account is required to make contributions, but users can choose to remain anonymous or use a pseudonym for the public data release.
- **Purpose Limitation & Licensing:** Users must be able to specify their preferences for how their data is used (e.g., for public display and evaluation vs. for model training). This is captured using (new) IETF AI Use Preferences and Creative Commons Preference Signals. We will discuss below how this might extend to other preference signal proposals and/or technical approaches to gating data.
 - This is critical for answering a likely FAQ around public AI data – if you succeed in creating actually useful training data or new benchmarks, won’t private labs just immediately use that dat as well?

2 The Serverless + Git MVP

2.1 Overview

Our initial MVP of the flywheel is a “Serverless + Git Platform” approach. It is meant to be a robust and scalable starting point for the data flywheel. It strictly separates the “write” (contribution) and “read” (display) components of the system.

The overall goal is to enable opt-in contributions of data (prompt/output + optional meta-data) with the state-of-the-art (caveat: also experimental / untested) preference signals and enforcement via per-item Creative Commons license (in MVP: default is CC0, CC-BY, CC-BY-SA) and an AI preference signal (using IETF aipref draft spec + CC Preference Signals draft spec). The data is distributed via:

- Hugging Face (gated) — full bundles by license bucket; access via request/terms; HF workflows used for deletions where possible.
- Public site — leaderboards + full dataset access; Cloudflare WAF/bot controls to discourage bulk scraping.

For full details, see the separate repo and its readme: #todo fill me in

For a short summary of the technical approach, see below bullet point:

- Frontend: A Next.js application hosted on Vercel, providing a simple, static interface for contributions.
 - No major lock-in here. Can easily be swapped for a lightweight static site, other modern web tech, etc.
- Authentication: User identity is managed via Auth.js (Next-Auth), using Hugging Face (HF) as the exclusive OAuth provider. This allows for clear attribution of contributions to a user’s public HF username (if they choose to).
- Data Storage: The single source of truth is a Hugging Face Dataset repository, which functions as a “Git-as-a-database.”
 - Starting with HF as it is a platform with specific focus on AI datasets and open culture
- Waiting room approach: The implemented workflow follows a two-stage “waiting room” pattern to ensure data quality and safety:

- How contribution works:
 - A user logs in with their Hugging Face account.
 - The interface accepts contributions in three ways: by uploading an OpenWebUI JSON export file, via URL params, or by pasting plain text.
 - The frontend parser automatically detects the OpenWebUI format, extracts meta-data (like `model` and `tags`), and prepends it as YAML frontmatter to the chat content.
 - Users choose a label for the chat type (“Good Chat” / “Fail Chat”) and attach (1) a Creative Commons license and (2) an IETF aipref and/or Creative Commons Preference Signal to signal preferences around AI use of the contribution.
 - A pseudo-anonymity system allows users to contribute publicly with their HF username, as “anonymous,” or with a custom pseudonym.
 - An informed consent checkbox, linked to Terms of Service and FAQ pages, is required for all submissions.
 - Upon submission, a serverless function writes the contribution not to the final dataset, but as a new, single JSON file in a `_waiting_room/` directory within the Hugging Face repository. This operation is fast and avoids write conflicts.
- How processing of the “waiting room” works:
 - A separate, asynchronous script is run on a regular schedule (e.g., as a daily GitHub Action).
 - This script fetches all pending files from the `_waiting_room`.
 - It validates each contribution and includes a placeholder for future content moderation and PII checks.
 - Validated contributions are batched and appended to a final, organized dataset file (e.g., `data/2025-08.jsonl`). Contributions are further bucketed by license/prefs.
 - To ensure atomicity, all file additions (to the final dataset) and deletions (from the waiting room) are performed in a single commit to the Hugging Face Hub. *#todo* when scaling, need to consider race conditions around the processing!

2.2 Advantages of a Serverless + Git approach

A serverless + Git stack keeps the “write path” lightweight for contributors and cheap to operate. Functions spin up on demand and idle to zero, so we can avoid paying for boxes that sit around; the trade-off is cold starts, which are well-documented and can be mitigated with provisioned concurrency when needed.

On the “read path,” a static site on a global CDN gives instant distribution and low operational overhead. Pages (e.g., from [Cloudflare](#), something like GitHub Pages, or similar) can read directly from the Git “source of truth” and serve assets from edge locations by default, which is exactly what we want for a browsable leaderboard and dataset browser.

Additionally, using the Hub (Git-backed) as the source of truth buys us a public audit trail and first-class versioning semantics. HF’s dataset repos are literally Git + LFS, with revision pinning via commit/tag/branch; storage is backed by object storage and scales. That maps cleanly to our “waiting room → release buckets” workflow and makes it easy to diff changes over time. (relevant HF docs: [datasets](#), [storage](#))

Moderation and PII handling are naturally centralized in the processing step. Because we trigger the write as a small file into a staging path and move it during a scheduled job, we can run filters, de-dup, and attach license/pref metadata before publication without asking contributors to learn tooling.

Basic safety and access controls are pragmatic at the edge. Cloudflare’s WAF/Bot Management and newer “AI bot” controls give us a reasonable anti-scraping posture for public downloads and pages, even if nothing on the open web is truly copy-proof. Recent product updates explicitly target AI crawlers, with default blocking and challenge flows we can enable. ([Cloudflare Docs](#), [WIRED](#), [Business Insider](#))

Finally, the “preferences and licenses” story fits the stack. Dataset cards and metadata natively expose a `license` field and other tags (Hub UI and YAML), and CC licenses give clear obligations (e.g., BY, BY-SA) we can enforce in packaging and docs. That lets us partition releases by license and publish compatibility notes in a way downstream users can actually follow. ([Hugging Face](#), [Hugging Face](#))

2.3 Disadvantages (and what to watch)

Trust centralizes in the processor. Contributors have to believe the middle layer won’t silently drop or reshape submissions. If/when we introduce event-driven ingestion (queues/streams), we must design for retries and duplicates.

UX isn’t perfectly “instant.” There’s an inherent gap between a user pressing “share” and seeing their item on the public site, because we run validation and batching. That’s a conscious choice, but we should set expectations and likely provide some kind status expectations.

Operationally, serverless isn’t “no ops,” it’s “different ops.” Cold starts exist (especially on sporadic paths), API limits are real on the platforms we hit (GitHub has documented ceilings; HF also rate-limits writes/reads even if specifics vary by endpoint), and “pay per use” can surprise us at scale without cost guardrails (see e.g. [GitHub Docs](#), [GitHub Docs](#)).

There’s also platform coupling to be consider. Using specific hosted CI/CD, serverless run-times, and hub APIs creates a degree of vendor lock-in; this is a known trade-off in serverless architectures and something we can blunt with portable formats (JSONL), documented exports, and “boring” interfaces (Git). ([CNCF](#))

Compliance remains non-zero. Deletions across mirrored artifacts (Hub revisions, static site snapshots, downstream forks) take a clear policy and a repeatable playbook. For licenses,

we're responsible for honoring CC obligations in our packaging and comms (e.g., keeping BY attribution fields intact; not mixing BY-SA content into incompatible bundles). CC's legalcode and guidance make those obligations explicit; our tooling should, too. ([Creative Commons](#))

Net: Serverless + Git gives us a pragmatic bridge—fast contributor UX, public versioning, cheap distribution—while we invest in moderation, idempotent processors, and clear license lanes. If we communicate the review gap, publish the processor's rules, and keep escape hatches (export scripts, mirrors), the trade-offs are acceptable for an MVP and refine-able over time.

2.4 Other reading:

- <https://arxiv.org/abs/2109.02846>
- <https://datascience.codata.org/articles/10.5334/dsj-2021-012>

3 Data Flow and Data Retention

Last updated: 2025-08-03

This document describes a data retention policy for both to Public AI Data Flywheel and the broader Public AI Chat frontend project. We describe the data we collect, when it's produced, how long we keep it, how contributor license & preference signals work (beta), and how our distribution channels operate.

To contextualize the data retention policy, it may be useful to first review a short bullet point, plain language description of all the ways that users create data in their use of the Public AI Chat Frontend, across 4 distinct websites:

- User visits the “landing page” `{{landing__page__link}}`
 - user can enter a query as guest
 - * created: a chat object
 - * chat is sent to model provider
 - * chat is stored, associated with generic guest account
- user can proceed to “the main app” `{{app__link}}`
 - user enters name, email, password. Handled by OWUI auth code.
 - user can enter queries. Creates a chat object associated with the user.
 - * see OWUI chat schema for all possible fields
 - * key fields: free text entry, feedback data
 - * chats are stored in the OWUI database so that users can browser their chat history, but are NOT used for training or eval.
 - * some standard data is stored for basic admin/engineering needs (request volume, errors, etc.), with short retention
- in current dev version, user can also share to openwebui.com if they wish to (requires an account with the openwebui community platform)
- in launch version, user can share to flywheel, on a separate <https://optinflywheel.com> (`#todo` exact url)
 - on sharing, a public chat object is created. This object will live in a gated HF repo and a public static site

4 Glossary of Defined Terms (for this Chapter)

“{The Public AI Chat Frontend}” or “Frontend” means the hosted interface described in this document that allows users to issue prompts to third-party model endpoints.

“Open WebUI Instance” or “OWUI” means the hosted Open WebUI application at `{{app_link}}` (or successor URLs) that provides optional accounts and chat history.

“Opt-in Data Flywheel” or “Flywheel” means the separate contribution and distribution platform at <https://optinflywheel.com> (or successor URLs) through which users may opt in to contribute data for public evaluation and research use.

“Provider(s)” or “Model Endpoint(s)” means third-party model services (e.g., national labs, commercial providers) that receive user prompts and return model outputs. The Frontend is a gateway to these services and does not control their retention, training, or use practices.

“Model Gateway” means the service layer that forwards requests from the Frontend to Providers and records a Request Envelope (metadata such as request ID, model ID, token counts, latency, and status).

“Request Envelope” means non-content request metadata retained for reliability, capacity, and SLO monitoring.

“Session Telemetry” means minimal first-party analytics collected on page load/navigation (e.g., timestamp, pseudonymous session ID, coarse locale, feature flags).

“Security Logs” means IP address and User-Agent records used for rate-limiting and abuse detection.

“Chat Object” means prompt(s), tool calls (if any), and model output(s) associated with a session or account within the Open WebUI Instance.

“Contribution” means any data a user intentionally submits to the Flywheel (e.g., prompt/output pairs, tags, corrections), along with per-item License and AI Preference Signal selections.

“AI Preference Signal” means an AI-use preference value the contributor attaches to a Contribution (e.g., IETF AI Preferences draft values and/or Creative Commons preference signals), intended to be conveyed downstream.

“License” means the Creative Commons license selected by the contributor for a Contribution (supported in the MVP: CC0-1.0, CC-BY-4.0, CC-BY-SA-4.0).

“License Bucket(s)” means the partitioning of Contributions into separate release artifacts by License (e.g., v1.0-cc0, v1.0-cc-by, v1.0-cc-by-sa).

“Waiting Room” means the Flywheel’s gated staging directory (e.g., `_waiting_room/` in a Hugging Face repository) where Contributions are first written prior to validation and release.

“Release” means a published version of the dataset (and associated notes/checksums) assembled from validated Contributions, partitioned by License.

“Gated Repository” means the Hugging Face dataset repository that requires an access request and acceptance of dataset-specific terms.

“Static Site” means the public site that hosts leaderboards and provides full dataset access, with anti-scraping controls.

“Anonymized Contributor ID” or “Pseudonym” means a non-identifying handle published with a Contribution when a contributor elects anonymity or a pseudonym instead of a Hugging Face username.

“Personal Data” means information that identifies or can reasonably be linked to an individual; “Sensitive Personal Data” means Personal Data that is sensitive by law or policy (e.g., health, financial, precise location, government identifiers).

“We/Us/Our” means [ENTITY NAME], the operator of the Frontend, the Open WebUI Instance, and the Flywheel.

“You/Your” means the individual using the Frontend and/or contributing to the Flywheel, or the entity on whose behalf the individual acts.

4.1 1) System scope & components

- Model Endpoints (National Labs/Providers): Third-party endpoints; {The Public AI Chat Frontend} is a gateway.
- Hosted Open WebUI: Browser interface for issuing prompts to supported models; optional user accounts.
- Opt-in Data Flywheel: Contributors may donate prompt/response pairs plus optional metadata for evaluation and research.
- Distribution Channels
 - Hugging Face (gated): Live view into dataset; access requires request and acceptance of terms.
 - Flywheel Static Site (public): Leaderboards and full dataset access for benchmarking; protected with Cloudflare anti-scraping controls.

4.2 2) What data is produced & when

4.2.1 2.1 Open WebUI (no account required, but optional and recommended)

- Session telemetry (minimal) #todo: Double check telemetry and provide a sample payload to show interested users;
 - Produced when: Page load / navigation
 - Includes: Timestamp, pseudonymous session ID, coarse locale, feature flags
 - Stored in: First-party analytics
 - Access: Eng/Analytics (aggregated)
 - Default retention: Raw 30 days; aggregates 13 months
- IP & User-Agent (security) #todo double check this
 - Produced when: Each request
 - Includes: IP, User-Agent for rate-limit/abuse detection
 - Stored in: Security log store
 - Access: SRE/Security
 - Default retention: 7 days, then delete/aggregate
- Query content
 - Produced when: On submit
 - Includes: Prompt + output
 - Stored in: interface database
 - Access: server admin only
 - Default retention: user can delete at any time; deleted along after 30 days of user inactivity. #todo, discuss this
- Error logs (sanitized) #todo double check this, provide example? least important to provide an example here.
 - Produced when: On error
 - Includes: Stack trace, request ID (no prompt/output bodies)
 - Stored in: Log store
 - Access: Eng (least privilege)
 - Default retention: 30 days

4.2.2 2.2 Open WebUI Accounts (optional)

- Profile and settings #todo double check this
 - Produced when: Sign-up
 - Includes: Email/OAuth ID, display name
 - Stored in: User DB
 - Access: Support/Eng
 - Retention: Kept until the user deletes it. If the account is inactive for 30 days, we delete the account and associated records.

4.2.3 2.3 Model Gateway (to providers)

- Request envelope
 - Produced when: Each request
 - Includes: Request ID, model ID, token counts, latency, status
 - Stored in: Metrics DB
 - Access: Eng/SRE
 - Retention: 90 days (aggregates 13 months)

Provider transparency: We do not guarantee provider behavior. We clearly display the exact payload we forward (headers + body summary) and link to the provider's own terms and policies where available. Users should review provider terms before use.

4.2.4 2.4 Opt-in Data Flywheel

- Contribution payload
 - Produced when: On explicit opt-in
 - Includes: Prompt, model output, optional tags
 - Stored in: Gated staging (effectively public) → license-bucketed release (also effectively public)
 - Access: Public
 - Retention: Indefinite
- HuggingFace username OR Anonymized contributor ID
 - Produced when: On ingest
 - Includes: user provides their huggingface username or selects a pseudonym (can leave blank)
 - Stored in: Dataset metadata

- Access: Public
 - Retention: Indefinite
 - Release artifacts
 - Produced when: On release
 - Includes: JSONL/TSV, checksums, notes
 - Stored in: HF (gated) and Static Site (public)
 - Access: Public (per channel)
 - Retention: Indefinite
-

4.3 3) Event timeline (how data flows)

1. User prompts a model → payload sent to model provider; query and response are saved if user made an optional OpenWebUI account.
 2. If user wishes to select a chat to share via opt-in, they can do so. They never have to. Licensing and preference signals are set at opt-in time.
 3. Opt-in chat goes to “waiting room”
 4. Processing script collects items from waiting room, applies some checks (PII, content moderation), moves them into release buckets.
 5. A separate processing script takes data from the gated HF repo and builds the static site with leaderboard and data dump
 6. Post-release: approved removals are honored in future releases and our mirrors; prior downloads may persist.
-

4.4 4) Licensing & Preference Signals (Beta)

For each contribution, the user selects:

- A Creative Commons license: CC0-1.0, CC-BY-4.0, or CC-BY-SA-4.0.
- An AI preference signal: an IETF AI preference (draft) value and/or CC preference signal (“IETF AI Pref combo”).

How we use these:

- We record `license` and `ai_pref` per record.
- Records are partitioned by license into separate release buckets.

- Downstream users must comply with the license; we publish compatibility notes (e.g., ShareAlike).
- Users may set account defaults; each submission can override.

Binding effect: Once a record is included in a release, that license applies to that copy.

4.5 5) Retention schedule

- Web edge
 - Data: IP + UA
 - Purpose: Abuse prevention
 - Retention: 7 days raw → delete/aggregate
 - Deletion path: Automatic purge
- Web UI
 - Data: Minimal telemetry
 - Purpose: Reliability metrics
 - Retention: 7 days raw → delete/aggregate
 - Deletion path: Automatic purge
- Gateway
 - Data: Request envelopes
 - Purpose: Capacity/SLOs
 - Retention: 90 days raw; 13 months aggregates
 - Deletion path: Automatic purge
- Open WebUI Accounts
 - Data: Profile and preferences
 - Purpose: Auth/consent
 - Retention: Kept until user deletes; 30-day inactivity → delete
 - Deletion path: Self-service delete / auto-purge
- Flywheel staging bucket (the “waiting room dir”, on HF)
 - Data: Pending contributions
 - Purpose: Moderation/de-duplication
 - Retention: Indefinite (moved to organized buckets in same repo)
 - Deletion path: Manual removal on request

- HF (gated)
 - Data: Released records (by license)
 - Purpose: Research access
 - Retention: Indefinite
 - Deletion path: Exclude from future versions; coordinate HF deletion where possible
 - Static Site (public)
 - Data: Leaderboards and full dataset access
 - Purpose: Benchmarking/transparency
 - Retention: Indefinite
 - Deletion path: Update/remove in future releases; past downloads may persist
-

4.6 6) Distribution & access control

4.6.1 6.1 Hugging Face (gated)

- Separate artifacts and checksums per license bucket and version (e.g., `v1.0-cc0`, `v1.0-cc-by`, `v1.0-cc-by-sa`).
- Access requires a request with acceptance of dataset-specific Terms of Use (no re-identification; honor license and AI preferences).
- Deletion requests: Where possible, we use Hugging Face-native workflows (e.g., issue/repo requests, maintainers' takedowns) to process deletions and exclude items from future versions.

4.6.2 6.2 Flywheel Static Site (public)

- Hosts leaderboards and provides full dataset access for building benchmarks.
 - Cloudflare anti-scraping posture: Bot Management/WAF rules, rate limits, Turnstile/JS challenges on download routes, tokenized short-lived URLs, robots/meta controls, pagination/throttling, and anomaly monitoring.
 - Reality check: Anti-scraping reduces but cannot prevent copying of public data. We limit risk via license partitioning, logged access flows, and clear terms.
-

4.7 7) Provider transparency (no guarantees)

- We do not control third-party model providers' retention or training practices and make no guarantees about them.
 - For every request, we show users a payload transparency panel (headers summary + body size/tokens) and, where possible, a link to provider terms.
 - Users choose whether to proceed with the selected provider in light of those terms.
-

4.8 8) Contributor rights & controls

- Opt-in only for contributions to the flywheel.
 - Per-item license and AI preference (beta).
 - Withdrawals: Submit a verified request to exclude your items from future releases. We'll also file/route requests through Hugging Face where applicable.
 - Account deletion & inactivity: If you keep an Open WebUI account, we keep your account data until you delete it; if you are inactive for 30 days, we delete the account and associated records. Released dataset copies remain under their chosen licenses.
 - Sensitive content: Do not submit personal/sensitive data; automated filters and moderation apply.
-

4.9 9) Roadmap

- HF-auth contributions: Let contributors submit via their Hugging Face accounts to preserve reputation and contributor stats.
 - AI preference UX: First-class UI for selecting license + IETF AI Pref combo, plus validator and compatibility helper.
 - Provider term links directory: Central index of provider policies; per-model tooltips in the WebUI.
-

4.10 10) Contacts

#todo

- Data removal / privacy: {{PRIVACY_EMAIL}}
- Security: {{SECURITY_EMAIL}}
- Research & benchmarks: {{RESEARCH_EMAIL}}

5 Rationale (and other flywheel variants)

5.1 Purpose of this section

This section explains the why behind {The Public AI Chat Frontend}'s data and distribution design, and lays out alternative governance paths and a future work (in particular, a focus on futures that involve healthy data markets, data intermediaries, federated learning, etc.)

We also discuss why we think an approach that includes a minimal retention frontend + opt-in flywheel platform can serve as a pragmatic bridge to future models: e.g., we can use the concepts used here to move towards independently governed data co-ops, eventual **federated learning, etc.

5.2 More on all the other approaches we could've taken

In choosing our architecture, we consider the following questions:

5.2.1 Where does the “final” data live?

- **Centralized Database:** Traditional server-controlled storage (PostgreSQL, MongoDB, etc.)
- **Public Repository:** Version-controlled platforms (GitHub, GitLab, Hugging Face Hub)
- **Distributed Network:** Peer-to-peer systems (IPFS, BitTorrent)
- **Totally Local:** Federated model where data stays on user devices
- **Something hybrid?:** Metadata centralized, actual data distributed

5.2.2 When is the user prompted to contribute?

- **Proactive:** User initiates contribution unprompted (e.g., “Share this chat” button)
- **Reactive:** System prompts based on signals (e.g., after thumbs down, ask “What went wrong?”)
- **Passive:** Automatic collection with prior consent (e.g., telemetry, browser extension)
- **Scheduled:** Regular prompts (e.g., weekly “best conversations” review)

- **Task-Based:** Specific requests for data types (e.g., “Help us improve math responses”)

5.2.3 What information object is created?

- **Simple Signal:** Binary feedback (/), star ratings, or flags
- **Annotated Conversation:** Full chat with user corrections, ratings, or notes
- **Preference Pair:** A/B comparisons between responses
- **Synthetic Example:** User-created prompts and ideal responses
- **Structured Feedback:** Form-based input (error type, severity, correction)
- **Multimodal Bundle:** Text + images + voice + metadata
- **More advanced structure data ...**

5.2.4 When is the data processed?

- **Pre-submission:** Client-side processing before data leaves user’s device
- **On-submission:** Real-time processing during the contribution flow
- **Post-submission:** Batch processing after data is received
- **Pre-publication:** Review and processing before making data public
- **On-demand:** Processing happens when data is accessed/downloaded

(In practice, there may be some processing at various steps, but it is important to clarify this to users)

5.2.5 How is the data accessed?

- **Direct Download:** Raw access to complete dataset (with rate limits)
- **API Access:** Programmatic access with authentication and quotas
- **Static Site:** Read-only web interface with anti-scraping measures
- **Gated Access:** Application/approval process for researchers
- **Hybrid Access:** Public samples + gated full access, or public metadata + restricted content
- **Streaming Access:** Real-time feeds for continuous model training

5.2.6 How much friction is acceptable?

- **Zero-Friction:** One-click actions with no interruption
- **Low-Friction:** Modal popup or inline form
- **Medium-Friction:** Redirect to separate interface
- **High-Friction:** Multi-step process, account creation, or technical skills required

5.3 Some Categories of Architectural Models

5.3.1 Standard “PrivateCo” Web App

An obvious option is to simply build a hosted “standard” “PrivateCo” /start-up style web app. In fact, in some contexts it may make sense to skip building an opt-in flyhweel and simply use the data generated by users directly for training, eval, etc. While one could argue that the Terms of Service for many existing tech products do make these products “opt in” in some sense, there are also serious downsides to the status quo (see e.g. [Fiesler, Lampe, and Bruckman 2016](#).)

While perhaps some users might prefer even prefer a start-up style model, we believe this would not be a good starting place for a public AI interface. We also believe it’s important to communicate to users how the public AI interface differs from e.g. using ChatGPT, Gemini, or AI overviews via search.

How this approach answers the above questions:

- **Where data lives:** Centralized database
- **When prompted:** Proactive (user initiates)
- **Information object:** Annotated conversations with structured feedback
- **When processed:** On-submission + pre-publication review
- **How accessed:** API + static site with rate limits
- **Friction level:** Medium (redirect to platform)
- **Pros:** Full control over UX, rich features, easy user management
- **Cons:** High maintenance, single point of failure, trust requirements
- **Example Stack:** Django/Rails + PostgreSQL, Next.js + MongoDB

5.3.2 Git/Wiki Platform

- **Where data lives:** Public repository
- **When prompted:** Proactive (user initiates)
- **Information object:** Markdown-formatted conversations
- **When processed:** Pre-submission (user does it) + CI/CD validation
- **How accessed:** Direct download via Git + web interface
- **Friction level:** High (technical knowledge required)
- **Pros:** Maximum transparency, built-in versioning, low cost
- **Cons:** Excludes non-technical users, limited data types
- **Example Stack:** GitHub/GitLab + CI/CD validation

5.3.3 Direct Telemetry

- **Where data lives:** Centralized analytics database
- **When prompted:** Passive (continuous collection)
- **Information object:** Simple signals with context IDs
- **When processed:** On-submission (real-time pipeline)
- **How accessed:** Aggregated dashboards only (no raw access)
- **Friction level:** Zero
- **Pros:** Massive scale, unbiased sampling, real-time insights
- **Cons:** Limited richness, privacy concerns, no corrections
- **Example Stack:** ClickHouse/BigQuery + streaming pipeline

5.3.4 Hybrid Model

- **Where data lives:** Centralized database
- **When prompted:** Reactive (triggered by signals)
- **Information object:** Signals + optional full conversations
- **When processed:** Signals processed immediately, conversations reviewed
- **How accessed:** Public aggregates + gated conversation access
- **Friction level:** Zero, then low
- **Pros:** Balances volume and quality, efficient targeting
- **Cons:** Complex implementation, two-system maintenance
- **Example Stack:** Telemetry backend + web app frontend

5.3.5 Serverless + Git Platform

- **Where data lives:** Public repository
- **When prompted:** Proactive or reactive
- **Information object:** Structured data files (JSON/YAML)
- **When processed:** On-submission via serverless function
- **How accessed:** Git access + static site generation
- **Friction level:** Low (automated complexity)
- **Pros:** Transparency + usability, serverless scaling
- **Cons:** Cold starts, API rate limits, complex error handling
- **Example Stack:** Vercel/Netlify + GitHub API + Hugging Face Hub

5.3.6 Federated Learning Model

- **Where data lives:** User devices (distributed)
- **When prompted:** Passive with consent

- **Information object:** Model gradients or aggregated statistics
- **When processed:** Pre-submission (on-device)
- **How accessed:** Only aggregated model updates available
- **Friction level:** Zero after setup
- **Pros:** Maximum privacy, no data transfer, infinite scale
- **Cons:** Complex implementation, limited debugging, device requirements
- **Example Stack:** Flower/TFF + edge deployment

5.3.7 Browser Extension (just a data ingestion choice: can be used with various backend choices above)

- **Where data lives:** Centralized or distributed
- **When prompted:** Proactive or passive
- **Information object:** DOM captures, interaction logs, selections
- **When processed:** Pre-submission (client-side) + server validation
- **How accessed:** Depends on storage choice
- **Friction level:** Low after installation
- **Pros:** Cross-platform, rich context, works anywhere
- **Cons:** Browser-specific development, installation barrier
- **Example Stack:** WebExtensions API + backend API

5.3.8 P2P Network Model

- **Where data lives:** Distributed across peer nodes
- **When prompted:** Passive (background sharing)
- **Information object:** Torrent-style data chunks
- **When processed:** Pre-submission by contributor + network validation
- **How accessed:** P2P client required for full access
- **Friction level:** Medium (client installation)
- **Pros:** No infrastructure costs, censorship resistant
- **Cons:** Availability issues, complex coordination
- **Example Stack:** libp2p + BitTorrent protocol + DHT

5.4 Terse decision Matrix for the above dimensions

Architecture	Data Location	Prompt Timing	Object Complexity	Processing Stage	Access Method	Friction	Best For
Web App	Centralized	Proactive	High	On-submission	API + Static	Medium	Full-featured contributions
Git/Wiki	Public repo	Proactive	Medium	Pre-submission	Direct Git	High	Technical community
Telemetry	Centralized	Passive	Low	Real-time	Aggregated only	Zero	Large-scale signals
Hybrid	Centralized	Reactive	Variable	Mixed	Tiered access	Variable	Balanced approach
Serverless + Git	Public repo	Proactive	High	On-submission	Git + Static	Low	Transparency + usability
Federated	Distributed	Passive	Low	On-device	Model updates only	Zero	Privacy-first
Extension	Variable	Variable	High	Client + server	Variable	Low	Cross-platform capture
P2P	Distributed	Passive	Medium	Pre + network	P2P client	Medium	Decentralized commons

5.5 Scenario Walkthroughs: A Practical Comparison

Here, we walk through two common scenarios and describe what happens (in one sentence) for each of the architectures described above.

#todo: these could be made crisper to highlight the key differences better (But also be honest about where there are similarities)

5.5.1 Scenario A: User marks a chat as “Good” – when does processing happen?

- **Web App:** Redirects to platform, PII scrubbed on submission, available via API after review
- **Git/Wiki:** User removes PII manually, creates PR, instantly visible on merge
- **Telemetry:** Signal sent, processed in real-time, only visible in aggregates
- **Hybrid:** Signal sent immediately, full chat processed if shared
- **Serverless+Git:** Modal appears, serverless function strips PII, PR created automatically
- **Federated:** Local processing only, contributes to next model update
- **Extension:** Captures state, removes PII client-side, sends to chosen backend
- **P2P:** Processes locally, shares with peers who validate before propagating

5.5.2 Scenario B: User corrects a factual error

- **Web App:** Editor interface, toxicity check on submission, published after human review
- **Git/Wiki:** User edits markdown, CI/CD checks format, visible immediately on merge
- **Telemetry:** Only captures “error” signal, no correction possible
- **Hybrid:** Error signal triggers correction UI, correction queued for review
- **Serverless+Git:** Inline correction, automated PII/toxicity checks, PR needs approval
- **Federated:** Correction processed locally, differential privacy applied
- **Extension:** Highlights error, pre-processes correction, sends to backend
- **P2P:** Broadcasts correction, network consensus before acceptance

5.5.3 Scenario C: Accessing the contributed data

- **Web App:** Researchers apply for API key, public sees samples on static site
- **Git/Wiki:** Anyone can clone repo, but rate-limited through CDN
- **Telemetry:** Only aggregated statistics available via public dashboard
- **Hybrid:** Public can see signals dashboard, researchers apply for conversation access
- **Serverless+Git:** Public (or gated) repo with all data, static site with search/filter
- **Federated:** No direct data access, only model checkpoints released
- **Extension:** Depends on backend choice, typically follows that model
- **P2P:** Must run client to access network, can specify data sharing preferences

5.6 Frontier approaches data cooperatives, federated learning, and more

In many cases, users may want to have data governed by community organizations (e.g., organized by domain/region/language) that hold rights and decide release cadence, licensing

defaults, and benefit policies.

We note that because our implementation is built on top of open-source software, communities can easily choose to deploy their own OpenWebUI instance and their own data flywheel and effectively operate entirely parallel, self-governed instances. If they also choose to share opt-in data via similar licensing and preference signal approaches, such datasets could be easily merged – but with fine-grained adjustments to precise details (e.g., slight modifications on retention, access, release cadence, content moderation, and so on.) Of course, data co-ops may choose to use quite different technical stacks. This approach is just one among many.

It may be possible to also move from an opt-in data flywheel approach to a federated learning-first approach. Here, model training occurs across user or institutional nodes; only gradients/updates (with privacy tech) are centralized. The dataset remains partitioned or local; central custodian minimized. This approach would:

- Reduces central data custody and breach surface
- Aligns with data-residency and institutional constraints
- Enables “learning from data that can’t leave”

But has some major downsides / existing barriers:

- Harder reproducibility and data auditability
- Complex privacy stack (secure aggregation, DP, client attestation)
- Benchmarking must be redesigned (federated eval)

This is a bigger leap, but we believe it’s important to begin to think about how the implementation of the Public AI Data Flywheels might support communities wishing to transition towards an FL approach.

One rough sketch might look like: * Build the MVP defined in Chapter 2 * Ship license + AI-preference metadata (MVP). * Maintain gated HF releases and public leaderboards/full data access. * Publish provider-payload transparency and link to provider terms (no guarantees). * Process deletions via HF mechanisms when possible; keep our mirrors in sync. * Phase 1 — Co-op pilots * Charter one or two community co-ops; define bylaws, scope, and release cadence. * Spin up many instances of interface + flywheel combos (can fork software directly, or use similar approaches) * Establish a concrete sharing / merging plan * And beyond! * Once several independent data communities, are operated, it might be possible to move from lightweight sharing and merging to more serious federation with technical guarantees. Perhaps this might start with federated evaluation and then move to federated training. Much more to do here, out of scope for this document.

6 Appendix 1: LLM Data Schemas

Here, we describe many variants of LLM data. This will be relevant for when we extend the flywheel to include more types of data, and especially shift towards promoting the sharing (via opt-in flywheels, but also via new market mechanisms) of richer “content data”.

- **Open Web / Crawls**

- **WARC/WAT/WET**

- * *WARC* (container for HTTP request/response records) — spec & overview: IIPC WARC 1.1; Library of Congress format note. ([IIPC Community Resources](#), [The Library of Congress](#))
 - * *WAT* (JSON metadata extracted from WARC) and *WET* (plain text extracted from HTML) — Common Crawl guides. ([Common Crawl](#), [Common Crawl](#))

- **C4 (Colossal Clean Crawled Corpus)** — TFDS catalog & generator code. Fields are essentially clean text segments with basic metadata. ([TensorFlow](#), [GitHub](#))

- **The Pile** (22-source, mixed corpus) — paper & HTML view. ([arXiv](#), [ar5iv](#))

- **Encyclopedic / Books**

- **Wikipedia XML dumps** (page/revision XML; SQL tables for links) — Meta-Wiki dump format; Wikipedia database download. ([Meta](#), [Wikipedia](#))

- **Project Gutenberg**

- * *Books*: plain text/HTML master formats; ePub/MOBI derived. ([Project Gutenberg](#))
 - * *Catalog schema*: daily RDF/XML (also CSV) for metadata; offline catalogs. ([Project Gutenberg](#))

- **Scientific / Legal**

- **arXiv** (Atom/OAI-PMH metadata; bulk & API) — OAI-PMH + API docs; bulk metadata page. ([info.arxiv.org](#), [info.arxiv.org](#), [info.arxiv.org](#))
 - **JATS XML** (journal article tag suite) — NISO standards; NLM JATS site. ([niso.org](#), [jats.nlm.nih.gov](#))

- **Code**
 - **BigCode** — **The Stack** / **The Stack v2** (source files + license/provenance metadata; dedup variants) — HF datasets, project docs, arXiv overview. ([Hugging Face](#), [Hugging Face](#), [BigCode](#), [arXiv](#))
- **Forums / Q&A / Social**
 - **Stack Exchange dumps** (XML: Posts, Users, Comments, Votes, etc.) — SE Meta/docs & Data Explorer. ([Meta Stack Exchange](#), [data.stackexchange.com](#))
 - **Reddit**
 - * *API JSON* schema — official API docs & help. ([Reddit](#), [Reddit Help](#))
 - * *Pushshift* (historical dumps; research dataset) — site & paper. ([pushshift.io](#), [arXiv](#))
- **Instruction / Conversations (Post-training SFT)**
 - **OpenAI-style chat schema** (role-tagged: `system|user|assistant`, plus tool calls) — API reference. ([OpenAI Platform](#))
 - **Alpaca** (JSON prompts/instructions/outputs) — Stanford post & repo; cleaned community set. ([crfm.stanford.edu](#), [GitHub](#), [GitHub](#))
 - **Databricks Dolly-15k** (human-written instruction/response pairs) — repo. ([GitHub](#))
 - **OpenAssistant OASST1** (message-tree conversations with roles) — HF dataset card. ([Hugging Face](#))
- **Preference / Feedback (RLHF & DPO)**
 - **HH-RLHF** (Anthropic helpful/harmless, JSONL pairs: `chosen` vs `rejected`) — dataset repo readme. ([GitHub](#))
 - **DPO format** (prompt + preferred vs dispreferred response) — DPO paper. ([arXiv](#))
- **Multimodal (for VLMs/ASR)**
 - **LAION-5B** / **Re-LAION-5B** (image-text pairs with CLIP scores; links) — LAION posts. ([laion.ai](#), [laion.ai](#))
 - **Whisper** (weakly-supervised ASR; audio → text pairs) — paper & blog. ([arXiv](#), [OpenAI](#))
 - **HowTo100M** (YouTube instructional video clips + narrations) — project page & paper. ([di.ens.fr](#), [arXiv](#))
- **Math-reasoning (often for post-training/eval)**

- **GSM8K** (grade-school word problems; JSON) — repo & HF dataset card. ([GitHub](#), [Hugging Face](#))
- **MATH** (competition problems with step-by-step solutions) — paper & HF. ([arXiv](#), [Hugging Face](#))

- **Common storage containers**

- **JSON Lines** / **NDJSON** — jsonlines.org; ndjson spec. ([jsonlines.org](#), [GitHub](#))
- **TFRecord** — TensorFlow tutorial. ([TensorFlow](#))
- **Apache Parquet** — project site. ([Apache Parquet](#))

#todo check all refs

7 Appendix 2 — Preference Signals for AI Data Use (CC signals + IETF AI Preferences)

#todo: improve the references here to specific lines of IETF draft and the CC Preference Signals FAQ

- **What CC signals are** A Creative Commons framework for *reciprocal* AI reuse: content stewards can allow specific machine uses if certain conditions are met (e.g., credit, contributions, openness). Overview & implementation notes. ([homepage](#), [implementation](#))
- **Four proposed CC signals (v0.1)**
 - **Credit (cc-cr)** — cite the dataset/collection; RAG-style outputs should link back when feasible.
 - **Credit + Direct Contribution (cc-cr-dc)** — proportional financial/in-kind support.
 - **Credit + Ecosystem Contribution (cc-cr-ec)** — contribute to broader commons.
 - **Credit + Open (cc-cr-op)** — release model/code/data to keep the chain open. Source (draft repo & posts). ([GitHub](#), [Creative Commons](#))
- **IETF AI Preferences (aipref) — the transport & vocabulary**
 - **Vocabulary:** a machine-readable set of *categories* (e.g., `ai-use`, `train-genai`) and *preferences* (`y` = grant, `n` = deny) with **exceptions**. Drafts. ([datatracker.ietf.org](#), [IETF](#), [IETF AI Preferences Working Group](#))
 - **Attachment:** how to convey these preferences via **HTTP Content-Usage** header and **robots.txt** extensions. Drafts. ([datatracker.ietf.org](#), [IETF](#))
 - **Structured Fields:** uses RFC-standardized HTTP structured field values. ([datatracker.ietf.org](#), [datatracker.ietf.org](#), [rfc-editor.org](#))
 - **Robots Exclusion Protocol** baseline. ([datatracker.ietf.org](#), [rfc-editor.org](#))
- **Putting them together (content-usage expression)**
 - Shape:
`<category>=<y|n>;exceptions=<cc-signal>`

Example in **robots.txt** (allow everything, but *AI use denied unless Credit*):

```
User-Agent: *  
Content-Usage: ai-use=n;exceptions=cc-cr  
Allow: /
```

Example **HTTP header** (deny *gen-AI training* unless *Credit + Ecosystem*):

```
Content-Usage: train-genai=n;exceptions=cc-cr-ec
```

(Syntax and examples from CC & IETF drafts.) ([Creative Commons](#), [IETF](#))

- **Operational notes (for this repo’s flywheel)**

- **Per-record fields** to store: `license` (CC0/CC-BY/CC-BY-SA) and `ai_pref` (IETF aipref value + optional CC signal), plus optional `attribution` handle. (Aligns with CC write-ups & IETF drafts.) ([Creative Commons](#), [datatracker.ietf.org](#))
- **Placement:**
 - * *Location-based* signals via **robots.txt** for site/paths. ([datatracker.ietf.org](#))
 - * *Unit-based* signals via **HTTP Content-Usage** on dataset files and API responses. ([datatracker.ietf.org](#))
- **Interoperability expectations:** signals are normative *preferences*; adherence relies on ecosystem norms (similar to robots.txt & CC license culture). ([Creative Commons](#))

- **Context & momentum**

- CC’s 2025 launch posts; IETF WG activity updates (e.g., IPTC note). ([Creative Commons](#), [Creative Commons](#), [IPTC](#))

#todo check all refs

8 Appendix 3: Example Legal Terms

Modeled after Mozilla Common Voice terms. #todo: closer comparison.

8.1 Opt-in Data Flywheel — Legal Terms (Draft)

Effective: [DATE]

Through the Opt-in Data Flywheel, you may contribute chats, corrections, and related materials to build openly accessible evaluation sets and datasets.

You may participate only if you agree to these Opt-in Data Flywheel Legal Terms (the “Terms”).

1. Eligibility

The Flywheel is open to individuals who are the age of majority in their jurisdiction, or to younger participants with verified parental/guardian consent and supervision. You must also comply with Our Community Guidelines/Acceptable Use Policy ([LINK]).

2. Your Contributions; Licensing; AI Preferences

- 2.1 Opt-in Only. Submitting to the Flywheel is purely voluntary and separate from using the Open WebUI Instance.

- 2.2 License Grant (per item). For each Contribution, you select one License (CC0-1.0, CC-BY-4.0, or CC-BY-SA-4.0). You grant Us a non-exclusive, worldwide right to publish, reproduce, modify (solely for formatting, moderation, and aggregation), distribute, and sublicense the Contribution under the selected License. Once included in a Release, that License applies to that copy of the Contribution.

- 2.3 AI Preference Signals. If you attach an AI Preference Signal, We will transmit and display it with the Contribution and document how Our systems interpret such signals. We cannot guarantee that downstream users or Providers will honor such signals.

- 2.4 Assurances. You represent and warrant that (a) you have the necessary rights to your Contributions; (b) your Contributions do not infringe third-party rights; (c) you will not include Sensitive Personal Data; and (d) you will comply with Our Acceptable Use Policy.

3. Accounts; Attribution; Pseudonymity

- 3.1 Auth. Contributions require authentication (e.g., Hugging Face OAuth).
 - 3.2 Attribution. You may choose to publish under your Hugging Face username, under a pseudonym, or as

“anonymous.” 3.3 Leaderboards. We may publish contribution metrics (counts, languages, tags) with your chosen public handle. We will not publish your email address. 4. Processing; Waiting Room; Release

4.1 Waiting Room. Submissions write to a staging directory. 4.2 Validation. We may run automated and human review for formatting, de-duplication, PII/safety checks, and License/AI-preference validation. 4.3 Release. Validated items are appended to License Buckets (e.g., vYYYY-MM) and published to a Gated Repository and mirrored to the Static Site. 5. Distribution; Access Control

5.1 Gated Repository. Access requires acceptance of dataset-specific terms (e.g., no re-identification; respect License and AI preferences). 5.2 Static Site. Public access includes anti-scraping measures (WAF/bot management, rate limits, tokenized URLs). Copying cannot be fully prevented; rely on License controls for downstream obligations. 6. Deletions & Takedowns

6.1 Future-Only Removal. Upon verified request, We will exclude the identified Contribution(s) from future Releases and update mirrors where feasible. Past Releases and third-party copies may persist. 6.2 Hugging Face Workflows. Where possible, We will route or honor takedowns via the Hugging Face repository’s native workflows. 7. Provider Transparency (No Guarantees)

We forward prompts to third-party Providers. We display a payload transparency panel and link to Provider terms when available. We do not control Provider retention, training, or other uses of data once sent to them. 8. Privacy; Retention

Retention and access for telemetry, envelopes, accounts, staging, Releases, and the Static Site are governed by the Data Retention & Contribution Policy (Section 3). That Policy is incorporated by reference. 9. Communications

By creating an account or requesting repository access, you may receive administrative emails (e.g., access decisions, policy updates). 10. Disclaimers; Limitation of Liability; Indemnity

THE FLYWHEEL AND RELEASES ARE PROVIDED “AS IS.” TO THE MAXIMUM EXTENT PERMITTED BY LAW, WE DISCLAIM ALL WARRANTIES (INCLUDING MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE, NON-INFRINGEMENT). WE WILL NOT BE LIABLE FOR INDIRECT, SPECIAL, INCIDENTAL, CONSEQUENTIAL, EXEMPLARY, OR PUNITIVE DAMAGES. Our aggregate liability under these Terms will not exceed USD \$500 (or the maximum permitted by law if lower). You agree to indemnify Us for third-party claims arising from your Contributions or breach of these Terms. 11. Updates

We may update these Terms by posting a new effective date. Continued use after the effective date constitutes acceptance. 12. Termination

We may suspend or terminate access at any time. Contributions included in prior Releases remain available under their Licenses. 13. Governing Law; Venue

These Terms are governed by the laws of [LAW & VENUE], without regard to conflict-of-laws rules. Exclusive venue lies in the courts of [VENUE].

8.2 Frontend Instance

Open WebUI Instance — Terms of Use (Draft)

Effective: [DATE]

These Terms govern your use of Our hosted Open WebUI Instance at {{app_link}} (or successor URLs). 1. Eligibility; Community Rules

The service is available to individuals who are the age of majority in their jurisdiction, or younger participants with verified parental/guardian consent and supervision. You must follow Our Community Guidelines/Acceptable Use Policy ([LINK]). 2. Accounts; Content

2.1 Accounts Optional. You may use OWUI without an account; certain features (history, settings, opt-in share flows) require an account. 2.2 Your Content. Prompts and outputs in your account are stored as Chat Objects to provide history and UX features. They are not used for training or evaluation by Us unless you explicitly opt in via the Flywheel. 2.3 Feedback Data. Thumbs, flags, and similar signals may be stored to improve product reliability and moderation and are handled per Section 3 (Retention Policy). 3. Provider Transparency

OWUI forwards your prompts to third-party Providers. We display a payload transparency panel and, where available, links to Provider terms. We do not control Provider retention, training, or other uses of your data. 4. Privacy; Retention; Security

Retention, deletion, and access controls for telemetry, Security Logs, Request Envelopes, account data, and error logs are governed by the Data Retention & Contribution Policy (Section 3), incorporated here by reference. 5. Sharing to the Flywheel

Sharing to the Flywheel is separate and requires explicit opt-in with per-item License and AI Preference Signal selections. See the Flywheel Terms. 6. Acceptable Use

You agree not to: (a) upload Sensitive Personal Data; (b) violate laws or third-party rights; (c) attempt to reverse engineer or abuse rate limits; (d) circumvent access controls; or (e) interfere with service integrity. 7. Communications

If you create an account, We may send administrative emails (e.g., login links, security alerts, policy updates). 8. Disclaimers; Limitation of Liability

OWUI IS PROVIDED “AS IS.” TO THE MAXIMUM EXTENT PERMITTED BY LAW, WE DISCLAIM ALL WARRANTIES (INCLUDING MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE, NON-INFRINGEMENT). WE WILL NOT BE LIABLE FOR INDIRECT, SPECIAL, INCIDENTAL, CONSEQUENTIAL, EXEMPLARY, OR PUNITIVE

DAMAGES. Our aggregate liability will not exceed USD \$500 (or the maximum permitted by law if lower). 9. Updates; Termination

We may update these Terms by posting a new effective date. Continued use after the effective date constitutes acceptance. We may suspend or terminate accounts for any reason, including AUP violations or security risk. 10. Governing Law; Venue

These Terms are governed by the laws of [LAW & VENUE], without regard to conflict-of-laws rules. Exclusive venue lies in the courts of [VENUE].