

# **SSIS Blueprint #1 v1.4 — Audio Pipeline (Polished + Clarified, No Other Changes)**

## **0) What we are building**

SSIS Audio Pipeline is an offline-first, CPU-only, resumable processing spine that turns raw audio into durable, versioned artifacts for downstream ML and product features—designed for environments where power loss is normal, not exceptional.

### **Primary artifacts produced**

- **AudioAsset:** canonical identity + source metadata
- **Normalized WAV:** canonical audio derivative (22050 Hz mono PCM WAV)
- **FeaturePack:** compressed HDF5 containing:
  - log-mel spectrogram
  - YAMNet embeddings via ONNX Runtime CPU
- **Segments:** speech/music/noise regions + derived silence
- **PreviewCandidate:** selected preview window with heuristic confidence + fallback mode
- **PipelineJob telemetry:** per-stage logs, metrics, error taxonomy, attempts
- **FeatureSpec registry:** lookup table linking short aliases → full feature specs

### **Non-negotiable design doctrine**

1. Idempotency at every stage (safe to rerun after crash/power loss)
2. Atomic publish for every artifact:
  - write temp → flush (best-effort fsync) → rename to final
3. Stage locks with stale-lock reclamation:
  - lock key prevents duplicate work after restarts
4. Versioned contracts from day one:
  - every artifact has schema\_id and version
5. Resilience tests in the spec:

- o kill mid-write, restart, verify correctness
- 

## 1) Locked decisions (v1.4)

### Orchestration & persistence

- Task orchestration: Huey
- Queue backend: SQLite (durable queue state offline)
- Metadata DB: SQLite (assets, jobs, locks, artifact index, feature specs)
- Default concurrency: 1 worker (safe baseline; configurable)

### Audio + ML

- Canonical audio derivative: 22050 Hz mono PCM WAV
- Embedding model: YAMNet via ONNX Runtime CPU
- Segmentation model: inaSpeechSegmenter (speech/music/noise)
- Silence: derived (gaps + optional VAD assist)
- Feature storage: HDF5 + gzip (single writer discipline)
- Preview selection: heuristic v1 (boundaries + energy variance + embedding variance), fallback to intro

### Artifact evolution & reprocessing

- Never overwrite artifacts. New configs/models → new FeatureSpec → new artifacts.
  - Feature configs are referenced by a short alias to avoid long identifiers.
- 

## 2) Scope and boundaries

### In scope (Blueprint #1)

- Ingest: local path + file upload
- Normalize: decode → canonical WAV
- Feature extraction: mel + YAMNet(ONNX) → HDF5

- Segmentation: inaSpeechSegmenter + silence derivation
- Preview selection: heuristic + fallback
- Observability: SQLite jobs + JSON logs + error taxonomy
- Resilience test harness (power cut simulation)

## Out of scope (Blueprint #1)

- Training custom models / weak supervision
  - CLAP/MERT semantic search
  - Transcription/diarization (optional later)
  - Distributed multi-node processing
  - UI/permissions system
- 

## 3) Repository layout (authoritative)

```

ssis/
  audio_pipeline/
    README.md
    pyproject.toml

specs/
  audio_asset.schema.json
  pipeline_job.schema.json
  feature_pack.schema.json
  segments.schema.json
  preview_candidate.schema.json
  ingest_local_request.schema.json
  preview_feedback.schema.json      # optional future contract

app/
  config.py
  db.py
  models.py          # AudioAsset, PipelineJob, StageLock, ArtifactIndex, FeatureSpec
  schemas.py
  huey_app.py
  orchestrator.py
  utils/
    hashing.py
    audio_meta.py
    atomic_io.py
    checkpoints.py
    paths.py

```

```
services/
  ingest_api/
    main.py
  worker_decode/
    run.py
  worker_features/
    run.py
  yamnet_onnx/
    yamnet.onnx
    yamnet.onnx.sha256
  worker_segments/
    run.py
  worker_preview/
    run.py
```

```
data/
  audio/{asset_id}/
  features/
  segments/
  preview/
  logs/
```

```
tests/
  contract/
  e2e/
  resilience/
```

---

## 4) Canonical paths and publish rules

### Canonical storage paths

- Original audio copy:
  - data/audio/{asset\_id}/original.<ext>
- Normalized canonical derivative:
  - data/audio/{asset\_id}/normalized.wav
- Feature packs (versioned by feature\_spec\_alias):
  - data/features/{asset\_id}.{feature\_spec\_alias}.h5
- Segments:
  - data/segments/{asset\_id}.segments.v1.json
- Preview:
  - data/preview/{asset\_id}.preview.v1.json

## Atomic publish rule (mandatory)

Every artifact write must:

1. write to temp path in same directory
2. flush + best-effort fsync
3. rename temp → final (the publish boundary)

## HDF5 special rule

- Never write directly to final .h5.
- Always write to .h5.tmp, then rename to .h5 only when fully valid.
- If a power cut occurs mid-write, the temp file may remain and must be ignored/cleaned on restart.

## Single-writer rule

Only worker\_features writes HDF5. All other stages write JSON only.

---

## 5) FeatureSpec: readable IDs + short aliases (with immutability)

### Why

Readable feature\_spec\_id strings can become long; we keep them human-auditable while using short aliases for filenames and locks.

### FeatureSpec fields

- feature\_spec\_id (human-readable, canonical string)
- feature\_spec\_alias (short stable alias) = first 12 chars of sha256(feature\_spec\_id)
- created\_at, notes

### Example

- feature\_spec\_id: mel64\_h10ms\_w25ms\_sr22050\_\_yamnet1024\_h0.5s\_onnx
- feature\_spec\_alias: a1b2c3d4e5f6

### Table:

### feature\_specs

- alias (pk)
- feature\_spec\_id
- created\_at
- notes

## Locked immutability rule (v1.4)

feature\_spec\_alias is frozen at first registration:

- If alias does not exist → insert (alias → feature\_spec\_id)
  - If alias exists and spec matches → no-op
  - If alias exists but spec differs → hard error FEATURE\_SPEC\_ALIAS\_COLLISION (and log job failure)
- 

## 6) Data contracts (spine contracts)

Every artifact includes:

- schema\_id (e.g., segments.v1)
- version (e.g., 1.0.0)
- asset\_id
- computed\_at
- recommended: pipeline\_version, worker\_version

### FeaturePack contract (v1.4 must include)

- feature\_spec\_alias
- feature\_spec\_id
- embedding\_backend: "onnxruntime"
- embedding\_model\_id: "yamnet"
- embedding\_model\_hash: "sha256:<...>"
- dataset pointers:
  - hdf5\_uri: data/features/{asset\_id}.{alias}.h5#/<dataset>

---

## 7) Database entities (minimum viable)

### Tables

1. audio\_assets
2. pipeline\_jobs
3. stage\_locks
4. artifact\_index (recommended)
5. feature\_specs

### Stage locks (v1.4)

Lock key:

- (asset\_id, stage, feature\_spec\_alias|null)

Rules:

- Acquire before stage begins
- Release at end
- Reclaim stale lock if older than TTL (e.g., 10 minutes)
- Log all lock acquisition/reclamation in pipeline\_jobs.metrics\_json

---

## 8) Pipeline stages (authoritative execution plan)

### Stage A — Ingest (FastAPI)

Input: local path or upload

Actions:

- copy file into SSIS-managed path
- compute sha256 hash
- extract best-effort metadata (duration/channels/sr)
- enforce idempotency via (owner\_entity\_id, content\_hash)

Outputs:

- AudioAsset row
- PipelineJob(stage="ingest")
- enqueue orchestrator start

Failure codes: INGEST\_FAILED, FILE\_NOT\_FOUND, HASH\_FAILED

---

## **Stage B — Decode + Normalize (worker\_decode)**

Input: AudioAsset.source\_uri

Goal: produce canonical wav

Outputs:

- data/audio/{asset\_id}/normalized.wav (atomic publish)
- PipelineJob(stage="decode")

Resilience:

- process in chunks
- checkpoint every ~60s processed audio
- restart resumes safely

Error codes: CODEC\_UNSUPPORTED, FILE\_CORRUPT, FILE\_TOO\_SHORT

---

## **Stage C — Feature extraction (worker\_features)**

Input: normalized wav

Goal: mel + embeddings saved in HDF5

FeatureSpec resolution

- Determine the active feature\_spec\_id for v1.4 (locked default).
- Compute feature\_spec\_alias and upsert into feature\_specs under immutability rule.

Outputs

- data/features/{asset\_id}.{alias}.h5 (atomic publish via .tmp)

- FeaturePack metadata record (DB and/or JSON)
- PipelineJob(stage="features")

Config (locked default FeatureSpec)

- Log-mel: 64 mel, 25ms window, 10ms hop, sr=22050
- YAMNet embeddings: 1024-D, hop=0.5s, onnxruntime CPU

Validation

- no NaN/Inf → else FEATURE\_NAN
- shape checks

Error codes: FEATURE\_NAN, MODEL\_OOM, FEATURE\_EXTRACTION\_FAILED, FEATURE\_SPEC\_ALIAS\_COLLISION

---

## Stage D — Segmentation (worker\_segments)

Input: normalized wav

Model: inaSpeechSegmenter → speech/music/noise

Silence: derived from gaps and optional VAD assist

Output

- data/segments/{asset\_id}.segments.v1.json (atomic publish)
- PipelineJob(stage="segments")

Confidence semantics (clarified)

- confidence is heuristic, not calibrated logits.
- Include confidence\_type: "heuristic\_v1".

Error codes: SEGMENTATION\_FAILED

---

## Stage E — Preview selection (worker\_preview)

Inputs

- normalized wav

- segments JSON
- embeddings from HDF5 for a specific feature\_spec\_alias

### **Locked FeatureSpec selection rule (v1.4)**

Preview reads embeddings from the FeaturePack matching the active feature\_spec\_alias, chosen deterministically:

1. If config SSIS\_ACTIVE\_FEATURE\_SPEC\_ALIAS is set, use it if that FeaturePack exists for this asset.
2. Else use the pipeline default feature\_spec\_alias (v1.4 shipped default).
3. If the selected FeaturePack does not exist, fail fast with FEATUREPACK\_MISSING and the orchestrator must (re)run features.

### Algorithm v1

1. identify candidate boundaries:
  - pause boundaries (>200ms low-energy)
  - segment boundaries in speech-heavy regions
2. generate candidates (default 60s)
3. score candidates:
  - energy variance (RMS variance)
  - embedding variance (diversity proxy)
4. select best above threshold else fallback:
  - intro = first non-silent window

### Output

- data/preview/{asset\_id}.preview.v1.json (atomic publish)
- PipelineJob(stage="preview")

Error codes: PREVIEW\_LOW\_CONF (not fatal; triggers fallback), FEATUREPACK\_MISSING

## **9) Orchestrator semantics (Huey + SQLite)**

### **Responsibilities**

- Determine required stages based on artifact existence + artifact index
- Acquire stage locks and dispatch tasks in order
- Enforce retries/backoff
- Log every attempt in pipeline\_jobs
- Dead-letter repeated failures

## **Retry policy (locked)**

- 3 attempts
- delays: 60s, 300s, 900s
- dead-letter after failures with error taxonomy

## **Idempotency rules (locked)**

- If final artifact exists, stage is skipped.
  - Temp artifacts are ignored and cleaned.
  - If a lock is held and not stale, stage is not duplicated.
- 

# **10) Observability & debuggability (must ship in v1.4)**

## **Required stage metrics**

- ingest: file size, hash time, format guess
- decode: output duration, chunk count, resample time
- features: inference time, mel/embedding shapes, NaN/Inf count, spec alias/id
- segments: segment count, class distribution, flip rate
- preview: candidate count, best score, fallback\_used, spec alias used

## **Pipeline-level health**

- end-to-end latency (ingest → preview ready)
- success rate

- backlog depth
- error breakdown by code

## Minimal error taxonomy (locked)

- CODEC\_UNSUPPORTED
  - FILE\_CORRUPT
  - FILE\_TOO\_SHORT
  - FEATURE\_NAN
  - MODEL\_OOM
  - SEGMENTATION\_FAILED
  - PREVIEW\_LOW\_CONF
  - FEATUREPACK\_MISSING
  - FEATURE\_SPEC\_ALIAS\_COLLISION
- 

## 11) Resilience tests (required, in blueprint)

### Contract tests

- validate API + artifact JSON against JSON schemas

### End-to-end smoke test

- ingest → orchestrator → decode → features → segments → preview

### Load-shedding simulation tests (must ship)

Test harness must:

- kill worker mid-write:
  - during decode
  - during HDF5 .tmp write
- restart pipeline and assert:

- no corrupt final artifacts
  - temp files ignored/cleaned
  - stale locks reclaimed
  - pipeline completes successfully
  - idempotency prevents duplication
- 

## 12) Re-extraction and versioning policy (explicit)

### When feature config changes

- Create a new feature\_spec\_id
- Generate new feature\_spec\_alias
- Produce new HDF5 artifact {asset\_id}.{alias}.h5
- Keep prior artifacts intact
- ArtifactIndex tracks which FeaturePacks exist per asset

No backfill is forced; backfill becomes an explicit batch job later.

---

## 13) Optional future contract: Preview feedback (for learning later)

Define early so later blueprints can plug in without redesign.

Contract: preview.feedback.v1 (optional)

- asset\_id
- preview window
- listen/skip metrics
- user actions
- timestamp

Blueprint #1 does not require collection.

---

## 14) Implementation plan (exact build order)

1. Contracts + DB primitives
  - schemas in /specs
  - DB models: AudioAsset, PipelineJob, StageLock, FeatureSpec, ArtifactIndex
  - utilities: atomic\_io, checkpoints, hashing, paths
2. Ingest API
  - local + upload
  - idempotency enforcement
  - enqueue orchestrator
3. Orchestrator (real)
  - stage planning + dispatch
  - locks + stale lock reclamation
  - retries + dead-letter
4. Decode worker
  - canonical wav
  - chunking + checkpointing
  - atomic publish
5. Features worker (Option B ONNX)
  - mel + embeddings via onnxruntime CPU
  - HDF5 .tmp atomic publish
  - FeatureSpec upsert with immutability enforcement
  - validation + telemetry
6. Segments worker
  - inaSpeechSegmenter + silence derivation
  - heuristic confidence + post-processing
7. Preview worker

- deterministic FeatureSpec selection rule
- candidate scoring + fallback

## 8. Resilience harness

- kill/restart tests + lock reclamation assertions

## 9. MVP acceptance

- offline CPU run
- deterministic artifacts
- safe restart after interruption
- job telemetry sufficient for debugging

This document is only Blueprint #1 — the SSIS Audio Pipeline module.

It's 1 out of the 3.

- Blueprint #1 (this one): Audio Pipeline — ingest → normalize → features (mel + YAMNet ONNX) → segments → preview + observability/resilience.
- Blueprint #2 (next): LLM Safety + Query Compiler — natural language → typed/validated command plans (no direct execution).
- Blueprint #3 (after): Real-time Scoring + Explainability — fast scoring + reason codes + feedback loops.

All three will later be linked into one complete SSIS system, but the doc I just wrote to you is specifically the module blueprint for #1.