# Building an Automated Vocal-Technique Analysis Pipeline

You can build this as a **pipeline**: (1) ingest audio → (2) extract acoustic features → (3) segment into voiced frames/notes → (4) compute interpretable scores (placement/weight/timbre) → (5) optionally train ML models to map features to human-labeled "vocal technique" ratings.

The key is: **tone placement / weight / timbre are not single physical variables**—they're *perceptual constructs*. So you'll need (a) clear operational definitions and (b) either rules + thresholds or a supervised model trained on expert labels.

# Complete System Overview

## 01

### Frontend Interface

Upload WAV/MP3 files, select audio type (spoken vs sung), and optionally choose prompt type (sustained vowel, reading passage, song verse) for analysis consistency.

## 02

### Audio Preprocessing

Resample to target rate (16k–44.1k), apply loudness normalization, light noise reduction, and trim silence regions to prepare clean audio input.

## 03

### Source Separation

For singing with accompaniment, run vocal isolation to extract clean vocal track for accurate analysis.

## 04

### Voice Activity Detection

Split audio into voiced and unvoiced segments to focus analysis on relevant vocal regions.

## 05

### Pitch Tracking & Segmentation

Extract F0 contours; segment spoken audio by syllable/phrase boundaries, sung audio by stable pitch regions (note-like segments).

## 06

### Feature Extraction

Compute acoustic features per segment: spectral, harmonic, formant, and cepstral measurements.

## 07

### Scoring & Inference

Apply rule-based algorithms or ML models to generate interpretable vocal technique scores from extracted features.

## 08

### Report Generation

Output comprehensive analysis: visualizations, numeric scores, confidence intervals, and contextual explanations.

# Acoustic Features That Correlate With Vocal Technique

## A) Timbre Analysis

Timbre is captured through **spectral shape** and **harmonic content**. Core measurements include:

- Spectral centroid (brightness)
- Spectral slope/rolloff
- HNR (harmonics-to-noise ratio)
- MFCCs (cepstral coefficients)
- Formants F1–F4 and bandwidths
- Inharmonicity (roughness)
- Jitter/Shimmer (micro-instability)

> 🗒 **Timbre scoring (0–100):** Brightness = normalized spectral centroid + rolloff. Breathiness = inverse(HNR) + noise energy above 3–5 kHz. Warmth = stronger low harmonics + steeper spectral slope.

## B) Vocal Weight

"Weight" relates to **source strength** (glottal closure and harmonic energy), not just perceived loudness.

**Key features:**

- Energy of harmonics (H1, H2) relative to higher harmonics
- Spectral tilt (harmonic dropoff rate)
- CPP (cepstral peak prominence)
- H1–H2, H1–A2, H1–A3 ratios
- Normalized SPL/loudness

> 🗒 **Weight scoring (0–100):** Heavier = higher CPP, flatter tilt, lower H1–H2 (more closure), stronger mid harmonics. Lighter = steeper tilt, higher H1–H2, lower CPP.

## C) Tone Placement

Placement is operationalized through **resonance/formant patterns** and **energy distribution**:

- Formant locations and vowel consistency
- Singer's formant (2.5–3.5 kHz)
- Spectral balance across bands
- Nasality proxy (anti-resonances)
- Open quotient/pressedness proxies

> 🗒 **Placement scoring (0–100 forwardness):** Higher 2.5–4k energy + stronger ring peak + stable formants = "forward." Dominant low-mid energy + reduced upper-mid ring + formant drift = "back."

**Critical note:** Placement varies with vowel and pitch. Your analyzer must compute it per vowel/segment, then average by context for accurate results.

# Generating Interpretable Values From Acoustic Features

## Option 1: Rule-Based Scoring

**Fastest to ship, ideal for MVP**

1. Select 6–12 core acoustic features
2. Normalize via z-score against reference dataset
3. Combine with weighted sums into indices
4. Output scores (0–100) with confidence metrics

Generate scores for: Timbre (Brightness, Breathiness, Warmth, Roughness), Weight (Light–Heavy, Pressed–Breathy), Placement (Forward–Back, Ring, Nasal).

*Works well for MVP but won't match expert perception in all cases.*

## Option 2: Supervised ML

**Best for "coach-like" ratings**

Requires labeled dataset where experts rate clips on placement, weight, and timbre descriptors. Include separate labels for spoken vs sung, vowels, and pitches.
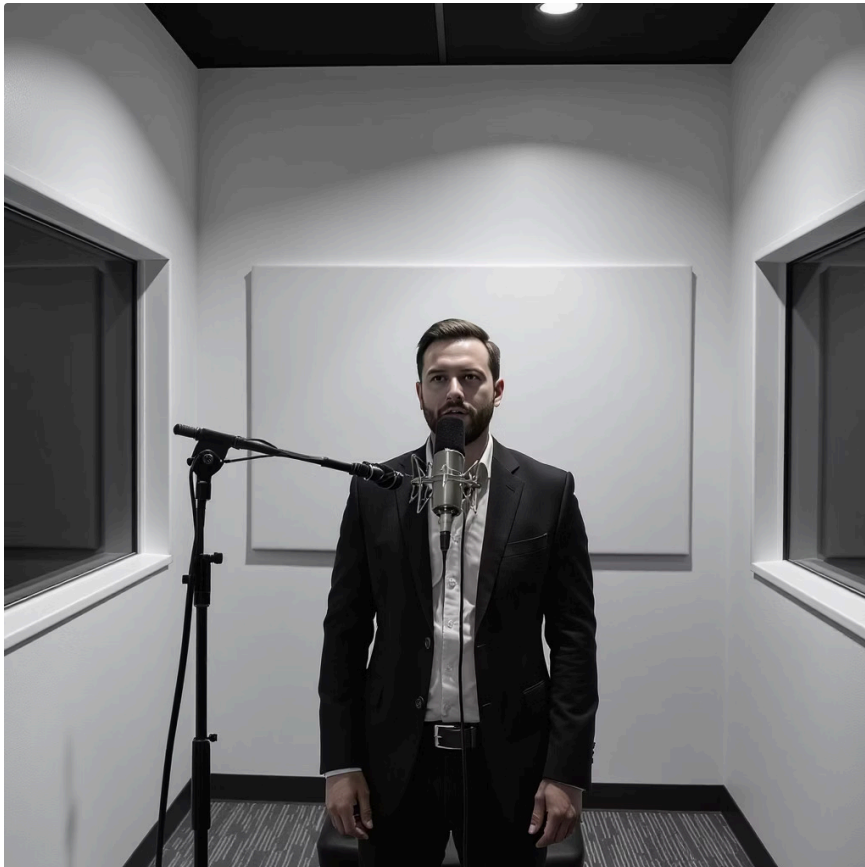
**Model architecture:**

- Inputs: log-mel spectrogram + F0 contour + summary features (MFCC, formants, CPP, HNR)
- Model: CNN/Transformer for embeddings + regression head
- Outputs: continuous scores (0–100) with uncertainty estimates

**Reality check:** Labeling is the hard part. Without consistent expert labels, model performance will be noisy.

# Spoken vs Sung: Critical Analysis Differences

Your analyzer must branch processing based on audio type to ensure accurate feature extraction and scoring. These contexts require fundamentally different analytical approaches.
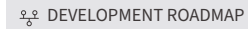


## Spoken Voice Analysis

- More consonants with less stable pitch contours
- Focus metrics: HNR, spectral tilt, formants per vowel, prosodic features
- Segment by vowel nuclei and voiced regions
- Account for natural pitch variation and speech rhythm
- Emphasize clarity and articulation measures

## Sung Voice Analysis

- Sustained vowels with stable F0 (optimal for analysis)
- Additional metrics: vibrato rate/depth, singer's formant/ring, note stability, register transitions
- Segment by note-like pitch regions
- Track consistency across sustained phonation
- Measure resonance optimization and harmonic richness

# Practical Build Plan: MVP to Production

## MVP Phase (2–4 weeks)

**Core functionality for rapid deployment:**

- Audio upload interface with basic UI controls
- VAD + pitch tracking implementation
- Extract MFCC, centroid, rolloff, HNR, formants, spectral tilt, CPP
- Generate three analysis panels: Timbre (brightness/breathiness/warmth), Weight (light-heavy, pressed-breathy), Placement (forward-back, ring index)
- Export PDF reports with visualizations

Delivers immediate value with rule-based scoring system.

**1**

**2**

## V2 Phase (Coach-Level Accuracy)

**Enhanced precision and intelligence:**

- Implement vocal isolation for mixed music tracks
- Advanced formant tracking (especially high notes)
- Dataset collection + expert labeling infrastructure
- Train regression models per task (spoken vs sung)
- Per-vowel/per-register analysis with confidence intervals
- Real-time feedback capabilities

Achieves professional-grade analysis matching expert human perception.

# Recommended Technology Stack

## Backend Architecture

**Primary:** Python with FastAPI for API endpoints

**Alternative:** Node.js + Python microservice for audio processing

Python provides superior DSP library ecosystem while FastAPI enables high-performance async request handling.

## DSP & Audio Processing

**Core libraries:**

- librosa – comprehensive audio analysis
- parselmouth – Praat integration for formants
- pyworld/crepe – F0 extraction
- torchaudio – PyTorch audio utilities

## Machine Learning

**Framework:** PyTorch for model training and inference

Use pretrained audio embeddings (wav2vec2-like models) as feature extractors, then fine-tune regression heads for vocal technique scoring tasks.

## Storage & Infrastructure

**Object storage:** S3-compatible service for features + reports

**Privacy-first:** Store extracted features and analysis results, not raw audio files (unless explicitly consented).

## Frontend Interface

**Framework:** React for component-based UI

**Visualization:** Waveform viewer + interactive plots (D3.js/Plotly) for spectral analysis, F0 contours, and score visualizations.

# Privacy & Consent: Non-Negotiable Requirements

If you let users upload voice recordings, robust privacy protections are legally and ethically mandatory. Voice data is biometric information subject to strict regulations.

### Explicit Consent & User Control

Require clear, informed consent for audio processing. Provide users with straightforward deletion controls and data access rights. Document exactly what happens to uploaded audio and how long it's retained.

### Encryption & Security

Encrypt all voice data at rest using industry-standard encryption (AES-256 minimum). Use TLS 1.3 for data in transit. Implement access controls and audit logging for all data access.

### Default Privacy Posture

Process audio, extract features, generate analysis report, then **immediately delete raw audio files**. Store only anonymized feature vectors and analysis results unless user explicitly opts into retention.

### Training Data Consent

If you plan to use user recordings for ML model training, make this **opt-in only** with separate explicit consent. Document training data usage in privacy policy. Provide compensation or incentives for contributed data.

**Regulatory considerations:** Depending on your jurisdiction and user base, you may need GDPR compliance (EU), CCPA compliance (California), HIPAA compliance (healthcare contexts), or other privacy frameworks. Consult with legal counsel early.

# Next Steps: Let's Get Specific

To provide you with maximally useful implementation details, I need one key decision from you:

# Simple Scorecard or Vocal Coach Tool?

Your answer determines the complexity, dataset requirements, and development timeline.

### If Scorecard Tool

I'll provide: Exact feature formulas with band ratios and thresholds, normalization procedures, JSON schema for output scores, rule-based scoring algorithms ready for implementation.

### If Vocal Coach Tool

I'll provide: ML model architecture specifications, dataset labeling guidelines, training pipeline design, suggested report layout with graphs and explanations, complete FastAPI endpoint design (upload → analyze → results).

Tell me your choice, and I'll deliver concrete, implementation-ready specifications tailored to your exact use case. We'll move from theory to working code.