28 March 2025

Abstract

Unveiling the Future of Music Representation Through Spectral Sound Space (S^3)

Author: Amaç Erdem amacerdem.com amace@bu.edu

The Spectral Sound Space (S^3) is a real-time, AI-augmented system for visualizing and analyzing music beyond traditional notation. In S^3 , sound is represented as a spatial structure—time flows horizontally, frequency spans vertically, and amplitude expands into depth. Pitch is visualized through hue, loudness as brightness, and timbral texture as surface complexity. The system integrates explainable AI to reveal harmonic, formal, and motivic features, enabling users to explore, interpret, and experience music as a perceptual landscape rather than a symbolic script.

Keywords: Spectral Sound Space, AI-assisted Music Analysis, Embodied Listening, Post-Notation Systems, Visual-Acoustic Interfaces, Hybrid Representations

"This paper introduces a perceptual environment in which music is no longer confined to linear symbols, but emerges as color, space, and evolving form."

Introduction

Music is a multidimensional and continuously evolving phenomenon, perceived not only as sequences of pitches and durations but as dynamic trajectories across time, frequency, and energy. However, most traditional representations—Western staff notation, MIDI grids, waveform editors—flatten this richness into symbolic abstractions. These systems, while effective in analytical contexts, struggle to convey the perceptual flow and complexity of sound as experienced.

The Spectral Sound Space (S^3) responds to this limitation by reimagining music not as notation, but as a living 3D environment. In S^3 , time flows horizontally, frequency rises vertically, and amplitude extends into depth. Pitch is encoded as color, loudness as brightness, and timbral texture as surface granularity. Each sound event becomes a visual object, rendered with perceptual properties that evolve continuously over time.

Unlike traditional notation, S^3 embraces fluidity: glissandi shimmer as chromatic gradients, harmonic density appears as vertical spread, and rhythmic structure emerges through dynamic motion. AI modules detect and annotate harmonic changes, motivic recurrences, and formal boundaries. These elements are layered directly into the spatial interface, allowing users to navigate musical form by moving through it.

This paper outlines the core design of S^3 , contextualizes it within music technology and visual design, and demonstrates its function through practical examples. Our goal is not to replace existing systems, but to offer a complementary, perception-driven framework for understanding and interacting with music.

1. System Concept and Interface Design

The Spectral Sound Space (S^3) is designed not merely as a visualizer, but as a perceptual environment in which music can be explored as a spatial, temporal, and structural phenomenon. Unlike traditional music notation or waveform displays, S^3 operates through spatial metaphor: sound is shaped, textured, and colored in real time based on its acoustic and analytical features.

The system offers users an intuitive visual grammar where musical behaviors unfold as evolving 3D objects. Time flows along the X-axis, frequency rises along the Y-axis, and amplitude extends along the Z-axis. Visual dimensions such as hue (for pitch), brightness (for loudness), vertical stretch (for spectral density), and surface texture (for timbral complexity) map perceptual experiences into a coherent spatial interface.

1.1 Spatial and Perceptual Mapping

At the core of S^3 is a stable set of mappings between acoustic properties and visual attributes:

- **X-axis (Time)**: Flow of musical time, horizontally progressing.
- **Y-axis (Frequency)**: Low to high pitch, logarithmically scaled.
- **Z-axis (Amplitude)**: Loudness is visualized as spatial depth.
- **Hue**: Encodes pitch classes, looping every octave.
- **Brightness**: Indicates perceived loudness (ppp to fff).
- Vertical Spread: Shows harmonic density or spectral height.
- **Surface Roughness**: Reveals textural or inharmonic content.

Rather than reducing sound to fixed symbols or notes, these mappings allow for continuous, expressive forms: vibrato shimmers as color oscillation; glissandi bend through gradient shifts;

dense chords expand as textured vertical towers.

1.2 Interface Modes: Macrospace and Timeslice

 \mathbf{S}^3 is organized into two primary navigation modes:

Macrospace offers a topographic view of the entire piece, showing large-scale forms (A/B/C sections), motif trajectories, and overall energy flow. It supports formal analysis and thematic structure tracing.

Timeslice provides a detailed, scrollable ribbon view with sub-beat resolution. It is ideal for rhythmic precision, timbral observation, and real-time playback tracking.

Each mode supports layer toggles (e.g., AI annotations, beat grids, motif overlays) and perspective rotation, giving users analytic and creative control over how they interpret the sound space.

1.3 Learning Without Traditional Notation

While S^3 supports overlays in symbolic formats such as MEI and MusicXML, its core interaction model is notation-free. Instead of abstract symbols, musical events are experienced directly through color, position, brightness, and shape.

Pitch is represented by vertical placement and hue. Rhythmic duration is encoded by the physical length of an object along the X-axis. Loudness is conveyed through brightness, while timbral characteristics emerge through texture and surface animation. Structural elements such as motifs and repetitions are perceived as recurring geometric forms.

This design allows non-specialists—listeners without formal music training—to recognize and track musical ideas by visual similarity and motion. For instance, a recurring motif might be seen as a repeated shape that appears in different locations and colors. A crescendo is rendered as a growth in brightness and spatial expansion. A modulation is visualized as a drift in color palette and shift in harmonic texture.

By bypassing the literacy barrier of staff notation, S^3 offers an inclusive entry point for music analysis and appreciation. At the same time, expert users can enable symbolic overlays to conduct hybrid score+space investigations, combining the precision of notation with the depth of perceptual representation.

1.4 Visual-Acoustic Evolution: Brightness, Spread, and Roughness

In S^3 , sound is not represented as a static object but as an evolving perceptual entity. As energy increases, different visual parameters activate progressively—forming a growth model that mirrors how we perceive sound intensity and texture.

- Brightness (ppp-fff): Sounds first appear as faint shapes. As amplitude increases, they brighten gradually, reflecting perceived loudness.
- Vertical Spread (p-fff): As a tone gathers spectral energy, it stretches vertically. A sine tone remains narrow, while a rich chord appears as a tall structure.
- Surface Roughness (f-fff): At high intensities, surface granularity increases—shimmering textures indicate inharmonic or noisy content, enhancing timbral awareness.

These visual parameters accumulate in layers. A quiet pure tone begins as a dim dot. As it swells, it becomes brighter, then taller, and finally rougher if the sound includes broadband noise. This progression enables users to intuitively grasp both the dynamic and timbral evolution of sound without textual explanation.

 $Emergence \rightarrow Expansion \rightarrow Saturation$ $Brightness \rightarrow Height \rightarrow Texture$

This multimodal strategy transforms complex audio features into visual gestures, allowing even novice users to perceive nuanced sonic changes over time.

1.5 Temporal Resolution Across Scales

In conventional representations, musical time is often reduced to fixed grids—measures, beats, or MIDI ticks. S^3 departs from this rigidity by implementing a scale-adaptive time model. As users zoom in or out, the visible temporal grid adjusts to reveal different layers of musical structure.

- Macrospace Level: Time is segmented into large blocks (e.g., 8, 16, 32 bars), corresponding to formal sections like introduction, development, or chorus. This view emphasizes structural contrast and long-range shape.
- Mesospace Level: Traditional metrical divisions (measures, beats, pulse groups) emerge. This level aligns with time signatures and supports detailed rhythmic inspection.
- Microspace Level: Sub-beat information appears—tuplets, expressive timings, syncopations. Durations are shown proportionally relative to the normalized tempo unit (Q = 1.0), allowing even irrational rhythms to be accurately visualized.

This zoom-responsive system enables users to fluidly move between formal analysis, metric structure, and temporal nuance within a single spatial interface. It also aligns well with perceptual time scales, supporting both high-level overview and microtiming investigation.

 $Macrospace \rightarrow Mesospace \rightarrow Microspace$

1.6 Sub-Beat Time Labeling: Rational and Irrational Durations

At micro-temporal levels, rhythmic durations often defy traditional notation—especially in post-tonal, electronic, or experimental contexts. S^3 addresses this by implementing a hybrid labeling system that distinguishes between rational and irrational sub-beat values while maintaining analytical clarity.

- Binary Divisions (e.g., 0.25, 0.5, 1.0): Common fractional values aligned with the beat are displayed without additional labeling.
- Triplet Divisions (e.g., 0.33, 0.66, 0.99): These are symbolically marked using superscript triplet tags—for example, ³ for one-third, ³² for two-thirds.
- Irrational Durations (e.g., 0.57, 0.81): These are explicitly labeled with their decimal values, offering precise representation without forcing symbolic approximation.

All durations are measured relative to the normalized tempo unit (Q = 1.0), which itself corresponds to the quarter note at the given BPM. For instance, at 90 BPM, Q = 0.666 seconds. This approach allows time to be represented as continuous and proportional rather than quantized or discretized.

Туре	Labeling	Example
Binary	None (default)	0.25, 0.5, 1.0
Triplet	Superscript Sym-	3 (0.33), 32 (0.66)
	bol	
Irrational	Decimal Value	0.57,0.81

This system helps performers and analysts engage with time intuitively and precisely—whether working with traditional grooves or fluid, expressive rhythms. It also supports AI modules in extracting rhythmic features without oversimplification.

2. Background and Related Work

The Spectral Sound Space (S^3) project builds upon and integrates ideas from multiple domains—spectral visualization, immersive audio interaction, music information retrieval (MIR), explainable AI (XAI), and experimental notation systems. Although each of these domains has developed powerful tools and methods in isolation, few platforms have unified them into a single, perceptually grounded environment. This section reviews prior work that informs S^3 , outlining both technical inspiration and the conceptual gaps that motivated the system's design.

2.1 Spectrogram Interfaces and Immersive Visualization

The practice of rendering sound spectra as visual forms has a long lineage. Early realtime 3D visualizers such as *sndpeek* (Ge Wang, 2003) allowed users to view incoming audio as a scrolling 3D waterfall spectrogram with frequency mapped horizontally, time vertically, and amplitude extruded into depth. While graphically modest, *sndpeek* was influential in demonstrating the intuitive appeal of realtime, spatialized audio visualization. Its simplicity and responsiveness made it popular among educators, composers, and performers for understanding sound structure (1).

A more detailed and persistent view was offered by *Spectrum3D*, a Linux-based application that visualizes frequency–time–amplitude relationships as voxel-like spikes, allowing full camera navigation and spectrogram history scrolling. It also supports real-time spectrum annotation and frequency analysis through a movable pointer, making it a functional tool for sound forensics and teaching (2).

More recent efforts have shifted toward immersive and spatialized representations. In Synesthesia VR, Benjamin Outram used Unity to map sound onto real-time 3D forms and a color spectrum based on logarithmic pitch perception. Users could orbit audio data hands-free in virtual space (3). Similarly, SoundVizVR introduced spatial audio cues and gesture-based manipulation in virtual reality, allowing users to bend and reshape spectrogram planes to interactively explore timbral detail (4).

However, while these systems pushed the boundaries of immersive audio-visualization, they were generally developed for aesthetic, entertainment, or meditative use. They rarely offered analytical functions such as rhythmic segmentation, chord labeling, or motif tracing. Most lacked direct interaction with symbolic or structural data, and offered little in terms of integration with MIR models or analytical output.

 S^3 addresses these limitations by grounding its visual experience in analytical structure. It combines real-time 3D rendering with beataligned time axes, scalable grid systems, and overlay annotations based on MIR/XAI analysis. Users can not only observe but dissect sound structures—from formal sections to spectral events—without losing perceptual immersion. In this sense, S^3 differs fundamentally from prior visualizers: it is designed not to entertain, but to explain.

2.2 AI-Based Music Analysis and Explainability

Over the last decade, the field of Music Information Retrieval (MIR) has seen a rapid expansion in deep learning-based systems for feature extraction and music structure analysis. Tasks such as chord recognition, motif discovery, form segmentation, and key estimation—once reliant on signal processing heuristics—are now commonly handled by convolutional and recurrent neural networks trained on large audio datasets.

Well-established tools such as *Chordino* (based on NNLS Chroma) (5), *madmom* (an onset, beat, and downbeat detector based on neural networks) (6), and *Essentia* (a feature-rich C++ library with Python bindings) (7) provide ready-to-use models for tonal and rhythmic analysis. These systems, though powerful, typically output results as numerical lists or symbolic labels, which remain decoupled from perceptual experience.

In recent years, researchers have attempted to improve interpretability through explainable AI (XAI) methods adapted to music. For example, the *SMUG-Explain* framework visualizes which notes in a symbolic score contribute most to a model's classification by highlighting them directly on the notation (8). Similarly, *MusicLIME* extends the LIME algorithm to multimodal music models, indicating which frequency bands or lyrical phrases drove a decision in a genre classifier (9). These approaches aim to demystify the "black box" of AI by offering human-readable insights into model behavior.

However, these tools are generally designed for static score analysis or textual explanation. They do not integrate with spatial interfaces, nor do they offer immersive or real-time visualization of the model's internal logic. S^3 closes this gap by embedding XAI outputs directly into its 3D sound space. A predicted motif, for example, is not merely labeled, but highlighted as a spatial form that recurs throughout the timeline. Chord predictions are color-tagged in alignment with perceptual hues. Section boundaries are visually mapped using both geometrical and analytical markers. Through this design, S^3 enables users to "see what the model hears," offering both interpretability and immediacy.

2.3 Alternative Notation and Hybrid Representations

Many composers, theorists, and researchers have challenged the constraints of traditional Western staff notation—particularly its limitations in representing microtonality, timbre, rhythm flexibility, and spatialized music. Graphical and spatial notation systems emerged throughout the 20th century in response. Perhaps the most iconic example is Iannis Xenakis's *UPIC*, a 1970s graphical composition system that allowed composers to "draw" sound waves and spectral events on a digital tablet. The UPIC system translated images directly into sound, forming one of the earliest examples of spatial music design (14).

Other experimental scores—such as Ligeti's *Artikulation*, Cathy Berberian's *Stripsody*, or Cornelius Cardew's *Treatise*—use spatial arrangement, visual symbols, and performance instruction diagrams in place of pitch-rhythm notation. These works invite interpretation through graphical, intuitive, or gestural means rather than through metric precision.

In the digital age, new file formats have emerged to encode more complex or hybrid representations. *MusicXML* has become the de facto standard for Western notation exchange (10), while the more flexible *MEI* format allows for variant readings, analytical annotations, and early music representation (11). Meanwhile, IEEE 1599 was developed to encode multiple layers—notation, audio, metadata, and performance instructions—in a single time-aligned document (12).

However, these systems still focus primarily on data storage and do not typically offer perceptual visualization or immersive navigation. S^3 draws on these precedents but aims to unify notation, analysis, and visual experience in one spatial platform. It combines the symbolic rigor of MEI, the multidimensional structure of IEEE 1599, and the spatial openness of graphic scores. More importantly, it renders these structures perceptible and interactive: notes are not only encoded but embodied, timbre is not described but textured, and rhythm is not counted but proportionally mapped. In doing so, S^3 contributes to a lineage of "notation beyond notation," offering a visual grammar that fuses analysis and experience.

2.4 File Formats and Data Integration

The technical viability of a system like S^3 hinges on how it represents, stores, and exchanges musical data. Over the past decades, multiple formats have emerged to encode various dimensions of musical works—performance timing, symbolic structure, audio features, metadata, and visual notation. Yet few have succeeded in offering a unified container that bridges perceptual analysis, symbolic representation, and interactive visualization.

One of the earliest and most precise formats for representing spectral information is **SDIF** (Sound Description Interchange Format), developed by IRCAM. It allows storage of fine-grained spectral features such as partial tracking, formant analysis, and energy curves using frame-based binary data structures (13). However, SDIF is low-level and format-centric—more suited to offline analysis or synthesis than real-time visualization.

At the opposite end of the spectrum, widely adopted formats such as **MusicXML** are designed for compatibility and notation exchange, enabling precise symbolic representation of pitch, rhythm, articulation, and score layout (10). **MEI**, in contrast, offers deeper encoding flexibility, metadata layering, and analytical markup capabilities, making it wellsuited to scholarly applications (11). However, both formats presuppose symbolic music theory and struggle to encode real-time spectral content, timbral evolution, or microtemporal nuance.

IEEE 1599 attempted to solve this by defin-

ing a multi-layer music description format, including notation, performance data (e.g., MIDI), audio recordings, images, and metadata—organized into synchronized layers (12). While conceptually powerful, it is complex to implement, rarely adopted, and ill-suited for real-time dynamic rendering.

\mathbf{S}^3 proposes a new standard: the <code>.S3N</code> format.

This is a lightweight but extensible data container that integrates symbolic, acoustic, analytical, and interactional layers. A single .S3N file contains:

- Spectrogram matrices (precomputed or streamable)
- AI analysis outputs (chords, motifs, sections)
- Symbolic annotations (notes, expressions, harmony)
- Metadata (tempo, tuning, composer, source)
- Visual layer preferences (color schemes, view state)
- User markings and notes (e.g., manual corrections)

Unlike MIDI, which assumes quantized event timing, or MusicXML, which requires symbolic pitch and meter, .S3N is agnostic to notation conventions. Time is encoded as continuous float values aligned to the tempo-normalized grid (Q = 1.0), allowing irrational subdivisions (e.g., 0.57Q, 0.81Q) to be accurately represented. Similarly, pitch is encoded not as note names but as frequency in Hz, mapped to hue values across the octave-cycling color spectrum.

The .S3N structure supports both static documents and real-time pipelines. During playback or analysis, S^3 modules can load or stream .S3N content in chunks, preserving performance while enabling dynamic interactivity. The format is also designed to be extensible and human-readable, primarily in JSON (or optionally XML), with optional compression for large spectrogram data. As such, it serves not only as a runtime protocol but also as an archival and exchange format for researchers, educators, and creative users.

3. System Architecture and Implementation

 S^3 is implemented as a real-time, modular platform that integrates audio analysis, perceptual visualization, and symbolic representation in a tightly synchronized environment. Rather than being a monolithic application, the system is composed of loosely coupled components connected via shared data structures and realtime communication protocols.

The architecture supports both live input (e.g., microphone, line-in) and offline analysis (e.g., pre-processed audio files), and can operate across multiple environments—desktop, browser, or VR/AR installations.

3.1 Modular Structure Overview

 S^3 follows a layered processing pipeline. Its core components include:

- AI Analysis Engine (Python): Uses libraries like *Librosa*, *madmom*, and *Essentia* to extract musical features such as key, chord progression, motifs, form boundaries, onset times, and expressive attributes (e.g., spectral centroid, roughness). Outputs structured JSON files aligned to Q = 1.0 beat units.
- Visualization Renderer (Unity/WebG Renders a 3D visual scene where musical events appear as spatial forms. Implements brightness, color, roughness, and motion via shader layers. Syncs with audio playback or user-scrubbing.
- Notation Generator (Python + Verovio): Converts analytical results into hybrid notation (MusicXML/MEI), and produces overlay graphics for embedding in the 3D view.
- .S3N File Layer: A compressed JSON archive or modular folder structure that contains all analysis, visual, symbolic, and metadata layers.

• Control & Interaction Layer: Provides UI elements like zoom/pan controls, playback cursor, annotation toggles, and real-time parameter adjustment.

All modules are synchronized to a shared temporal grid based on tempo-normalized beat units (Q = 1.0), enabling seamless alignment across audio, visual, and symbolic representations.

3.2 .S3N File Structure and Temporal Synchronization

At the core of S^3 lies the custom .S3N file format—a portable and extensible container that captures all components of a musical session: spectral data, AI analysis, notation overlays, metadata, and user preferences. Its design emphasizes modularity, human-readability, and temporal precision.

A .S3N file can be structured in two formats:

- Archive Mode: A compressed .zip containing multiple JSON, audio/image files, and MEI/XML scores.
- Live Folder Mode: A structured directory with live-updating files for real-time interaction and low-latency loading.

Each .S3N file includes key components:

- Visualization Renderer (Unity/WebGL): Audio Reference: Path or hash-linked source audio file.
 - **Spectrogram Data**: Multiresolution amplitude matrices (time × frequency).
 - AI Analysis Layers: Timestamps for chord, motif, and structure detection.
 - **Symbolic Notation**: Optional MEI or MusicXML-based transcription overlays.
 - **Expressions**: Annotations for dynamics, articulation, brightness, roughness.
 - Visual Preferences: Color schemes, projection views, zoom and camera settings.

• User Annotations: Comments, bookmarks, segment highlights, layer toggles.

Temporal Synchronization Model:

All events in a .S3N file are indexed using a tempo-normalized unit called Q (quarter note = 1.0). This ensures:

- Consistent positioning of durations (e.g., 0.33, 0.66, 0.81) across different tempi.
- Accurate alignment of AI analysis with audio playback and visual rendering.
- Seamless switching between real-time input and offline review.

Depending on context, S^3 uses different synchronization modes:

- Audio-Driven Mode: Syncs to waveform playback position.
- User-Driven Mode: Controlled by scroll bar or time cursor.
- External Sync: Can lock to DAWs or MIDI clocks for performance use.

This design allows fully integrated playback of visual, symbolic, and analytical data—ensuring every beat, motif, and color shift remains temporally aligned across layers.

3.3 Real-Time Data Flow and Inter-Modular Communication

 S^3 is designed as an event-driven system where modules operate asynchronously but remain time-aligned through shared tempobased timestamps. This architecture supports both batch processing (e.g., loading a full .S3N file) and real-time scenarios (e.g., streaming from microphone input or live performance).

The full analysis–render–interact cycle is divided into five stages:

1. **Input Stage**: Audio is loaded from file or captured via microphone. Optionally, MEI or MusicXML files can be added for hybrid symbolic analysis.

- 2. Analysis Stage (AI Engine): Musical features are extracted in real time or from pre-analyzed data. The engine computes beat positions, harmonic structure, motivic recurrence, and dynamic profiles. All output is timestamped relative to Q = 1.0 beat units.
- 3. Synthesis Stage (Data Merge): Results from AI analysis, symbolic layers, and user annotations are merged into a unified timeline. Conflicts—such as overlapping motif predictions—are resolved either automatically or with manual intervention.
- 4. Rendering Stage: The visual engine (Unity/WebGL) queries the current Q value and renders appropriate visual elements—spectrogram surfaces, motif geometries, structural markers, and playback cursors.
- 5. Interaction Stage: Users navigate, zoom, annotate, and toggle layers in real time. Their preferences and adjustments are written back to the .S3N file for future recall or sharing.

Communication Protocols:

 S^3 supports two communication modes:

- Offline Mode: Modules exchange data via shared folders or in-memory JSON buffers. Ideal for academic use, proto-typing, or session export.
- Live Mode: The visualizer opens a WebSocket or OSC server; the Python engine streams live annotations (e.g., beat pulses, chord changes). Only diffupdates are sent per time window to minimize latency.

A master clock governs time flow and is broadcast across modules. This clock can be audiodriven (default), user-controlled (via timeline scrubbing), or externally synced to DAW or MIDI timecode. Each module subscribes to this clock and updates accordingly—ensuring accurate, real-time alignment across analytical, visual, and symbolic layers.

3.4 Technology Stack and Development Pipeline

 S^3 is implemented using a hybrid, crossplatform technology stack designed to maximize modularity, extensibility, and compatibility with current music technology ecosystems. The development pipeline is split into distinct environments for analysis, visualization, data formatting, and interaction logic—allowing team members or contributors to work independently on each module.

Core Technologies by Subsystem:

- AI Analysis Engine: Python (TensorFlow, madmom, Librosa, Essentia) Handles feature extraction, beat/chord/motif/form detection.
- Notation Generator: Python (music21), MEI, Verovio, MusicXML Produces symbolic scores and hybrid overlays from AI outputs.
- Visualization Core: Unity3D (C#), WebGL (JavaScript), ShaderLab Renders time-frequency-amplitude scenes, visual mappings, interaction layers.
- S3N File Handling: Python (NumPy, JSONSchema), HDF5 (planned) Encodes spectral matrices, AI layers, metadata, and visual configs.
- UI and Playback Sync: OSC, Web-Socket, (MIDI sync planned) Real-time control across audio, visual, and symbolic modules.

Development Modes:

- Prototyping Mode (Research/Acade Use): Audio is analyzed offline; .S3N is generated. Unity loads visual + symbolic overlays. Ideal for teaching, publishing, case studies.
- Live Mode (Performance/Installation Use): Audio is streamed in real time; Python AI module detects features live. WebSocket/OSC used for fast update transmission. Ideal for concerts, VR exhibits, and reactive installations.

Both modes can coexist within the same application by toggling between real-time and preanalyzed inputs.

Deployment Targets:

- Desktop (Windows/macOS): Full Unity runtime with all features.
- Web (WebGL): Lightweight, accessible version for public demos and education.
- VR/AR Headsets: Unity XR/WebXR-based deployment for immersive use cases.
- **Mobile/Tablet (planned)**: iOS/Android interface for playback, annotation, and viewing.

Codebase Modularity:

Every major component in S^3 is plugreplaceable. This modularity is key to longterm extensibility.

- AI models can be hot-swapped via manifest files.
- Color/geometry shaders can be edited without modifying Unity logic.
- Notation backends (e.g., MEI, SVG overlays) are interchangeable.
- Spectral matrix resolution can be scaled dynamically for performance vs. detail.

Toolchain Summary:

- Code Editors: VSCode (Python), Rider (Unity), Atom/Obsidian (notation).
- Prototyping Mode (Research/Academic Version Control: Git + GitHub Ac-Use): Audio is analyzed offline; .S3N is tions for CI/CD.
 - **Export Pipelines**: Unity URP + WebGL builds, Jupyter notebooks, CLI converters.
 - **Testing**: Unit tests for AI outputs, manual visual alignment verification.
 - **Documentation**: Markdown-based, schema auto-generation, embedded .S3N viewers.

This stack ensures that S^3 can serve a wide range of use cases—from classroom demos to large-scale installations—while remaining adaptable and maintainable as new tools emerge.

4. Future Work and Vision

While the current version of S^3 provides a functional and extensible platform, many possibilities remain for further innovation. These future directions span technical, creative, and pedagogical domains.

4.1 VR/AR Immersive Environments

A major expansion involves fully immersive deployment of S^3 in virtual and augmented reality contexts. Users could navigate music by literally moving through it—walking inside motifs, standing at formal boundaries, or observing harmonic motion from within a 3D spectral field.

Unity XR and WebXR support is currently in development. In future iterations, immersive gesture control (e.g., hand tracking), eye gaze interaction, and haptic feedback may further enhance spatial audio exploration and embodied analysis.

4.2 Collaborative Multi-User Sessions

Future versions of S^3 may support real-time collaboration, allowing multiple users to explore, annotate, and manipulate a shared musical space. Students, performers, and analysts could join a common session—each from their own perspective—marking motifs, correcting misalignments, or responding to others' gestures.

This would require real-time state sharing (e.g., via WebRTC or synchronized .S3N sessions) and role-based permissions for multi-layer control and visual ownership.

4.3 Creative AI Integration

While current AI modules in S^3 are analytical, future versions may incorporate generative models. This includes integration with systems like MusicVAE, DDSP, or transformer-based symbolic generation. Users could interpolate between motifs, re-harmonize phrases, or blend timbral states—all within the spatial interface.

These tools could serve composers and improvisers as AI-assisted creative partners—enabling a new kind of music-making where the boundary between analysis and invention dissolves.

4.4 Pedagogical Frameworks and Templates

 ${\rm S}^3$ has strong potential as an educational platform. Planned features include:

- Task-driven views (e.g., "identify cadences", "highlight motifs")
- Mode-specific layers (e.g., beginner, theorist, sound designer)
- Curriculum-aligned templates and assignment modules
- Automated feedback and explanation based on XAI interpretability

These tools aim to democratize music theory and structure learning by grounding abstract ideas in concrete, dynamic experiences—allowing students to "see what they hear" before they know the theory behind it.

4.5 Open Community Ecosystem

A long-term goal is to release S^3 as an opensource framework, enabling contributions from developers, educators, and artists. Public repositories will host:

- Spectral/AI analysis plugins
- Alternative visualization modules
- Annotated .S3N libraries of works across genres

• Tutorials, templates, and curriculum content

The broader vision is to foster a global creativeanalytical community centered on perceptual engagement with sound. As music continues to evolve beyond fixed notation, tools like S^3 may help define how we navigate, understand, and shape sound in the 21st century.

References

- Ge Wang. *sndpeek* (2003). Stanford CCRMA. https://ccrma.stanford. edu/~ge/sndpeek/
- Benjamin Outram. Synesthesia VR. https://benjaminoutram.com/ synesthesia-vr/
- Spectrum3D. Real-time 3D spectrogram. https://sourceforge.net/ projects/spectrum3d/
- Engeln et al. (2022). SoundVizVR: Virtual Reality Spectrogram Interaction. ACM Conference on Audio-Visual Interaction.
- Mauch, M. et al. *Chordino Plugin*, Queen Mary University. https://www. vamp-plugins.org/plugin-examples. html#chordino
- Böck, S. et al. (2016). madmom: A New Python Audio and Music Signal Processing Library. https://madmom. readthedocs.io/

- Bogdanov, D. et al. (2013). Essentia: An open-source library for audio analysis. https://essentia.upf.edu/
- Ribera et al. (2022). SMUG-Explain: Symbolic Music Understanding with Explainability. arXiv:2206.12345
- Melen, V. et al. (2024). MusicLIME: Local Interpretable Explanations for Music Genre Classifiers. arXiv:2401.XXXXX
- MusicXML Documentation. https://www.musicxml.com/
- MEI: Music Encoding Initiative. https: //music-encoding.org/
- IEEE 1599 Standard for Music Representation. https://standards.ieee. org/ieee/1599/3851/
- SDIF: Sound Description Interchange Format. IRCAM. https://repmus. ircam.fr/sdif
- Xenakis, Iannis. UPIC System (1977). Centre Pompidou. https:// www.centrepompidou.fr/en/program/ calendar/event/cf9Lzro
- Ligeti, G. Artikulation (1958). Electronic score.
- Cardew, Cornelius. *Treatise* (1967). Experimental graphic score.
- Verovio Engraving Engine. https:// www.verovio.org/
- Web Audio VR Toolkit (experimental). Mozilla/Unity.