INFER: Learning Implicit Neural Frequency Response Fields for Confined Car Cabin

Anonymous authorsPaper under double-blind review

000

001

002003004

010 011

012

013

014

016

017

018

019

021

024

025

026

027

028

029

031

032

035

037

038

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Accurate modeling of spatial acoustics is critical for immersive and intelligible audio in confined, resonant environments such as car cabins. Current tuning methods are manual, hardware-intensive, and static, failing to account for frequency selective behaviors and dynamic changes like passenger presence or seat adjustments. To address this issue, we propose INFER (Implicit Neural Frequency Response fields), a frequency-domain neural framework that is jointly conditioned on source and receiver positions, orientations to directly learn complex-valued frequency response fields inside confined, resonant environments like car cabins. We introduce three key innovations over current neural acoustic modeling methods: (1) novel end-to-end frequency-domain forward model that directly learns the frequency response field and frequency-specific attenuation in 3D space; (2) perceptual and hardware-aware spectral supervision that emphasizes critical auditory frequency bands and deemphasizes unstable crossover regions; and (3) a physics-based Kramers–Kronig consistency constraint that regularizes frequencydependent attenuation and delay. We evaluate our method over real-world data collected in multiple car cabins. Our approach significantly outperforms time- and hybrid-domain baselines on both simulated and real-world automotive datasets, cutting average magnitude and phase reconstruction errors by over 39% and 51%, respectively. INFER sets a new state-of-the-art for neural acoustic modeling in automotive spaces.

1 Introduction

Accurate modeling of acoustic environments is fundamental to diverse applications, including architectural design, immersive audio rendering, and human-computer interaction Koyama et al. (2025). While techniques for room acoustics and open-field settings are relatively mature, car cabins have recently emerged as a critical yet underexplored application space. Unlike conventional rooms, cabins are compact, irregularly shaped enclosures with a heterogeneous mix of reflective and absorptive materials, whose acoustic responses are further complicated by their highly dynamic usage—seats recline, windows open, passengers Yoshimura et al. (2012). These properties create transfer characteristics that are difficult to capture using traditional measurement or simulation pipelines. At the same time, the acoustic environment inside the cabin is becoming central to the in-vehicle experience, enabling high-fidelity entertainment and safety-critical spatial alerts. Modern audio pipelines, such as Dolby Atmos and Sony 360 Reality Audio, aspire to deliver immersive sound in vehicles but require precise characterization of these transfer functions. Existing approaches rely on labor-intensive manual tuning, extensive in-situ measurements, or costly simulations based on idealized CAD models, all of which degrade under real-world perturbations. These factors collectively motivate the need for a data-driven, physically grounded, and adaptive modeling framework that can generalize across diverse cabin conditions while preserving perceptual fidelity and spatial audio quality.

Recent advances in neural implicit representations (INRs) Molaei et al. (2023); Zhang et al. (2025) provide a compelling alternative to handcrafted acoustic models. INRs learn continuous, resolution-agnostic mappings from spatial coordinates to signal values using multilayer perceptrons, enabling geometry-free reconstruction of complex fields. Extensions such as Neural Acoustic Fields (NAF) Luo et al. (2022), INRAS Su et al. (2022), and AV-NeRF Liang et al. (2023) have shown

that impulse responses can be compactly encoded by learning emitter—receiver transfer functions directly from data. Yet these models predominantly operate in the time domain and treat individual frequency components uniformly. For car cabins, however, both physical acoustics and human perception demand frequency-selective modeling: low-frequency modes dominate room-scale resonances, speaker crossovers introduce mid-band artifacts, and perceptual salience varies with spectral weighting of phase and magnitude.

Learning acoustic frequency response fields directly in the frequency domain can lead to fine-grained and physically-aware modeling of confined, resonant spaces such as car cabins. By predicting each frequency bin independently, our approach naturally captures sharp spectral features and modal resonances that are often blurred in time-domain formulations. The spectral formulation also facilitates hardware-aware supervision: unreliable or unstable frequency bands can be identified and downweighted, and perceptually important regions—such as phase-sensitive low frequencies Oxenham (2018)—can be emphasized using auditory-inspired weighting schemes. Moreover, frequency-domain forward modeling encodes propagation delays exactly as phase shifts, eliminating the interpolation artifacts and discretization errors that plague time-domain approaches.

Recent methods have explored hybrid pipelines that uses both time-domain and frequency-domain forward modeling. For instance, AVR Lan et al. (2024) uses both sample delay and phase correction to model time-of-flight. However, AVR ultimately predicts time-domain impulse responses and assumes frequency-independent attenuation, limiting its ability to capture spectral variability across the scene. Furthermore, AVR employs uniform weighting across all frequencies during supervision, ignoring known variations in human perception and hardware characteristics.

We propose **INFER** (Implicit Neural Frequency Response fields), the first implicit neural representation framework that learns a frequency-domain acoustic field inside confined environments using a fully spectral modeling pipeline. Our formulation couples a differentiable frequency-domain renderer with a complex-valued neural network that predicts the frequency response field conditioned on emitter location, receiver location, and their direction. We propose a novel frequency-specific spectral weighting for both phase and amplitude, enabling perceptual and hardware-aware loss design. Finally, we introduce physically grounded consistency constraints derived from Kramers–Kronig relations to regularize the joint behavior of attenuation and phase delay, leading to more interpretable and physically plausible reconstructions. Our approach delivers consistent performance boosts across the spectrum, outperforming the closest baseline by 39% in magnitude error and 51% in phase error on average (Table 2). These gains reflect both spectral fidelity and directional accuracy, and are visually evident in Fig. 1.

Our key contributions in this paper can be summarized as:

- A novel end-to-end frequency-domain modeling framework that directly learns the frequency response field and frequency-specific attenuation in 3D space.
- A perceptually motivated, frequency-weighted supervision strategy that emphasizes critical bands and accounts for hardware artifacts such as crossovers and directivity lobes.
- A physically consistent formulation using Kramers–Kronig relations that jointly regularizes spectral amplitude and phase, enhancing generalization and interpretability.
- A comprehensive evaluation on both simulated and real car cabin datasets, demonstrating state-of-the-art phase and magnitude field reconstruction.

2 RELATED WORK

Neural Implicit Representations for Physical Fields. Neural implicit representations (INRs) have emerged as a powerful framework for modeling continuous physical signals by learning coordinate-to-signal mappings using multilayer perceptrons. Foundational methods such as SIREN Sitzmann et al. (2020) and Fourier feature encodings Tancik et al. (2020) allow compact modeling of high-frequency functions, enabling applications in 3D vision Mildenhall et al. (2021); Martel et al. (2021) and implicit surface reconstruction Wang et al. (2021). These ideas have been extended to domains like fluids Holl et al. (2020), mmWave Takawale & Roy (2025), and wave propagation Orekondy et al. (2023), demonstrating the versatility of coordinate-based learning. Our work

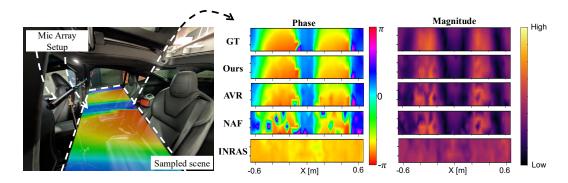


Figure 1: Acoustic field modeling inside a car cabin. Left: Measurement setup in the backseat of a Tesla Model X, where a speaker emits sound and spatial responses are recorded over a 2D grid. Middle: Spatial distribution of phase at 720Hz for ground truth (GT), our method (Ours), and baselines (AVR, NAF, INRAS). Right: Corresponding log-magnitude (energy) plots. Our method reconstructs smoother and physically consistent fields that preserve wavefront geometry and acoustic shadowing, outperforming baselines that exhibit artifacts or spatial inconsistency.

builds on these insights and targets learning frequency response fields of spatially varying acoustic transfer functions within a confined car cabin, where modal resonances, material absorption, and directionality jointly shape the acoustic field. This setting requires jointly modeling amplitude and phase attenuation over a frequency spectrum.

Neural Modeling of Acoustic Fields. Accurate modeling of sound fields in complex environments has traditionally relied on geometric acoustics, ray tracing, or numerical methods such as FEM and BEM Feng et al. (2013); Mo et al. (2015); Etgen & O'Brien (2007). These methods are physically grounded but computationally expensive and require detailed knowledge of geometry and material parameters. To address this, recent works like NAF Luo et al. (2022) and INRAS Su et al. (2022) have proposed neural implicit models to learn audio impulse responses from spatial data, often in time domain and with limited physical constraints. AVR Lan et al. (2024) introduces a hybrid frequency-time modeling pipeline but still lacks frequency-specific attenuation modeling. In contrast, our method INFER directly predicts frequency response fields, allowing frequency-aware supervision, perceptual spectral weighting, and physically grounded modeling of dispersion and absorption using Kramers–Kronig constraints—capabilities that existing impulse-response field prediction methods do not offer.

Car Cabin Acoustic Modeling and Applications. Acoustic field modeling inside car cabins presents unique challenges due to the confined space, material heterogeneity, complex modal behavior, and intricate reflection patterns Accardo et al. (2018); Yoshimura et al. (2012). Classical simulation techniques based on FEM or BEM Wang et al. (2013); Liu et al. (2017) are accurate but computationally prohibitive for design iteration or personalization. Empirical IR measurements Farina et al. (1998) and equalization techniques often ignore the global structure of the acoustic field, focusing instead on specific locations. Recent learning-based methods lack frequency-aware modeling and typically neglect physically grounded constraints essential for accurate modeling. Our method addresses these limitations by learning a continuous, physically grounded frequency-domain representation of the cabin's 3D acoustic field, enabling accurate reconstruction of both amplitude and phase, with explicit modeling of dispersion, crossover behavior, and material absorption—critical features not captured by prior empirical or neural approaches.

3 Primer: Physics of Acoustic Propagation in Lossy Media

3.1 PROBLEM PREMISE

Achieving physically consistent and perceptually accurate acoustic modeling in confined spaces like car cabins requires a deeper understanding of how sound propagates in complex, lossy media. Unlike free-field environments, car interiors exhibit complex modal behavior, intricate multipath interference, and frequency-dependent absorption—making frequency-domain analysis not just convenient

but essential. As established in Sec. 1, our method adopts a frequency-by-frequency modeling approach that gives the neural network the flexibility to understand these phenomena. To motivate and ground our spectral formulation, this section introduces the three key physical concepts underpinning our design: (1) how free-field propagation naturally maps to a frequency-domain formulation, (2) how rich multipath effects can be modeled via the Huygens–Fresnel principle, and (3) how signals interact in real media and experience attenuation and dispersion.

3.2 Free-Field Propagation and its Frequency domain Representation

To understand acoustic propagation, we begin with the classical time-domain free-space model. When a point source at $\mathbf{p}_{tx} \in \mathbb{R}^3$ emits an impulse at t = 0, the pressure at a receiver located at $\mathbf{p} \in \mathbb{R}^3$ in an ideal, lossless medium experiences decay in energy and arrives with a delay and is given by the 3D Green's function Kuttruff (2016):

$$h(t) = \frac{1}{4\pi r} \delta\left(t - \frac{r}{v}\right), \quad r = \|\mathbf{p} - \mathbf{p}_{tx}\|, \quad v = \text{speed of sound}$$
 (1)

The energy decay is due to 1/r spherical spreading of the pressure wave and the delay corresponds to the time-of-flight r/v. To arrive at the frequency domain representation of this phenomena, we apply the Fourier transform which yields:

$$H(f) = \frac{1}{4\pi r} \exp\left(-j\frac{2\pi fr}{v}\right) \tag{2}$$

The magnitude remains governed by 1/r, while the propagation delay is now expressed as a frequency specific phase shift $e^{-j\omega\tau}$, where $\omega=2\pi f$. This forms the building block of our frequency domain rendering approach.

3.3 THE HUYGENS-FRESNEL PRINCIPLE FOR MULTIPATH EFFECT MODELING

Acoustic wavefields in confined spaces arise from intricate multi-path interactions involving reflections, diffractions, and scattering. To capture this behavior, we draw inspiration from the Huygens–Fresnel principle Lian (2023), which posits that each point on a wavefront acts as a secondary emitter. In the frequency domain, the resulting complex pressure at a point \mathbf{x} can be modeled as:

$$P(\mathbf{x}, \omega) = \int_{\Omega} G(\mathbf{x}, \mathbf{x}'), S(\mathbf{x}', \omega), d\mathbf{x}', \tag{3}$$

where $G(\mathbf{x}, \mathbf{x}')$ is the Green's function encoding phase and amplitude propagation from \mathbf{x}' to \mathbf{x} , and $S(\mathbf{x}', \omega)$ is the frequency-domain strength of secondary emission. This formulation motivates our design: instead of tracing discrete reflection paths, we model the volume as a continuous field of directional secondary emitters. Each voxel learns to re-radiate incoming energy in all directions, capturing reverberation and scattering in a physically grounded, data-driven manner.

3.4 ATTENUATION AND DISPERSION IN MEDIA

Real acoustic environments are inherently lossy. As the sound propagates, amplitudes decay due to absorption and scattering (attenuation) and phases evolve at frequency-dependent speeds (dispersion). Crucially, these two effects are not independent artifacts - attenuation and dispersion are inherently linked. In any linear, time-invariant medium, the way amplitude varies with frequency determines how phase varies with frequency (and vice versa); one cannot be chosen independently of the other. This coupling is formalized by the Kramers–Kronig (KK) relation ODONNELL et al. (1981). Practically, this matters because past models fit only amplitude decay which, by construction, miss the paired frequency-dependent phase response that a real medium must exhibit.

Kramers–Kronig relations. The Kramers–Kronig relations express that the phase-bearing and amplitude-bearing parts of the medium's correction to wavenumber as Hilbert-transform. Physically, they ensure that no component of the response can occur before its excitation, i.e., *causality*. In acoustics, frequency-dependent propagation is written via a complex wavenumber $k(\omega) = k_0(\omega) + \delta(\omega)$, $k_0(\omega) = \omega/v$, where $\delta(\omega) = \text{Re}\left[\delta(\omega)\right] + j \text{Im}\left[\delta(\omega)\right]$ captures medium-induced modifications. The KK relations impose

$$\operatorname{Re}\left[\delta(\omega)\right] = \frac{1}{\pi} \mathcal{P} \int_{-\infty}^{\infty} \frac{\operatorname{Im}\left[\delta(\omega')\right]}{\omega' - \omega} d\omega', \quad \operatorname{Im}\left[\delta(\omega)\right] = -\frac{1}{\pi} \mathcal{P} \int_{-\infty}^{\infty} \frac{\operatorname{Re}\left[\delta(\omega')\right]}{\omega' - \omega} d\omega'. \tag{4}$$

Complex attenuation fields. We operationalize this principle by predicting, at each spatial point, a *complex attenuation* that separates amplitude loss and phase modification: $\delta(f, \mathbf{x}) = \sigma(f, \mathbf{x}) + j \beta(f, \mathbf{x})$, where $\sigma \ge 0$ is the absorption coefficient and β encodes dispersion-induced phase-velocity deviation. The KK-consistency is maintained through the Kramers–Kronig consistency regularizer explained in 4.3.

Physically consistent volume rendering. Once δ is known locally, its effects *accumulate* along a path as multiplicative transmittance and additive phase. For a small segment of length Δu ,

$$T_{\text{mat}}(\Delta u) = \exp(-\delta \, \Delta u) = \underbrace{\exp(-\sigma \, \Delta u)}_{\text{amplitude decay}} \cdot \underbrace{\exp(j \, \beta \, \Delta u)}_{\text{phase shift}}. \tag{5}$$

Over a full path $\mathbf{p}(s)$ of length L, amplitude and phase accumulate as $T_{\rm amp} = \exp(-\int_0^L \sigma \big(f,\mathbf{p}(s)\big)\,ds)$, and $\phi_{\rm mat} = -\int_0^L \beta \big(f,\mathbf{p}(s)\big)\,ds$. Prior acoustic neural fields typically model absorption or overall amplitude but ignore the causally paired, frequency-dependent phase response. In contrast, we are the first to encode KK-consistent complex attenuation in a neural acoustic renderer, preventing non-physical phase behavior and improving both interpretability and generalization.

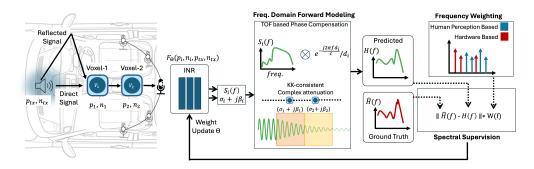


Figure 2: **System Overview.** Illustration of our frequency-domain acoustic forward model. For each point sampled along rays cast from the microphone, the MLP predicts a frequency-domain signal and attenuation. A TOF-based phase shift is applied to the signal and material based absorption and phase shifts are applied to produce the final response by accumulating signal from all directions.

4 METHODS FOR IMPLICIT NEURAL FREQUENCY RESPONSE FIELD

We introduce INFER, a fully frequency-domain neural rendering framework for modeling spatially varying frequency response fields in real-world environments. Unlike prior time-domain methods, INFER directly learns complex frequency responses—capturing sub-sample propagation delays, frequency-dependent attenuation, and dispersive phase shifts. This allows flexible perceptual weighting across frequencies, such as down-weighting hardware-specific crossover bands or emphasizing phase at low frequencies for localization. Fig. 2 gives a brief overview of INFER.

4.1 NEURAL FIELD PARAMETERIZATION

Given a scene with a sound source located at \mathbf{p}_{tx} and oriented along $\hat{\mathbf{n}}_{tx}$, we define a neural field F_{Θ} that predicts the local frequency-domain behavior at any spatial query point \mathbf{p} and frequency f:

$$F_{\Theta}: (\mathbf{p}, \hat{\mathbf{n}}, \mathbf{p}_{tx}, \hat{\mathbf{n}}_{tx}) \mapsto \{\delta(f, \mathbf{p}), S(f, \mathbf{p}, \hat{\mathbf{n}})\}$$
 (6)

Here, $\delta(f,\mathbf{p}) \in \mathbb{C}$ is the complex attenuation encoding the frequency-dependent transmission loss at point \mathbf{p} , and $S(f,\mathbf{p},\hat{\mathbf{n}}) \in \mathbb{C}$ is the directional spectrum re-radiated from that point toward direction $\hat{\mathbf{n}}$. Together, they fully characterize how an incoming acoustic wave is transformed and retransmitted from each location in the volume. Unlike other neural acoustic rendering models, INFER predicts all quantities directly in the frequency domain. The goal of the neural field is thus to answer: given a source at $(\mathbf{p}_{tx}, \hat{\mathbf{n}}_{tx})$, what frequency domain signal is re-emitted in any direction $\hat{\mathbf{n}}$ from point \mathbf{p} , and what is the frequency-specific material-induced attenuation along the way? We implement F_{Θ}

using a two-branch architecture. The first branch takes $(\mathbf{p}, \mathbf{p}_{tx})$ and predicts $\delta(f, \mathbf{p})$ via a material sub-network. The second branch conditions on the learned features from δ , the receiver direction $\hat{\mathbf{n}}$, and source direction $\hat{\mathbf{n}}_{tx}$, and predicts the directional retransmission spectrum $S(f, \mathbf{p}, \hat{\mathbf{n}})$. This decomposition reflects the physical structure of the problem: attenuation is direction-independent, while retransmission is highly directional.

4.2 Frequency-Domain Rendering

Our goal is to predict frequency response at any receiver location. To synthesize the acoustic frequency response at a receiver location \mathbf{p}_{rx} , we cast rays in direction $\hat{\mathbf{n}}$ and sample N points along the ray as $\mathbf{p}_k = \mathbf{p}_{rx} + u_k \hat{\mathbf{n}}$. At each sampled point, we query the neural field to evaluate local frequency-domain properties and accumulate their contributions using a physically motivated rendering equation:

$$H_{\hat{\mathbf{n}}}(f) = \sum_{k=1}^{N} S_k(f) \cdot \frac{1}{4\pi u_k} \cdot e^{-j2\pi f u_k/v} \cdot e^{j\phi_k(f)} \cdot \alpha_k T_k \tag{7}$$

This equation models how sound emitted from the transmitter propagates through the environment and contributes to the received frequency response along direction $\hat{\mathbf{n}}$. At each sampled point, $S_k(f)$ denotes the local directional spectrum predicted by the neural field, $e^{-j2\pi f u_k/v}$ introduces the phase shift due to time-of-flight delay, and $\frac{1}{u_k}$ accounts for spherical spreading via distance-based amplitude decay. The term $\alpha_k = 1 - \exp(-\sigma_k \Delta u_k)$ represents the discrete opacity arising from local absorption, while $T_k = \prod_{j < k} (1 - \alpha_j)$ captures accumulated transmittance from earlier samples along the ray. Finally, $\phi_k(f) = \sum_{j < k} \beta_j \Delta u_j$ models the cumulative phase shift induced by dispersive propagation through the medium. Together, these terms account for direction-dependent emission, distance-based decay, frequency-selective absorption, and phase dispersion—without discretizing time or relying on post-hoc transforms. To model a realistic microphone, which integrates sound from multiple directions, we perform weighted integration over a discrete set of directions: $H(f) = \sum_m G(\hat{\mathbf{n}}_m) H_{\hat{\mathbf{n}}_m}(f)$, where $G(\hat{\mathbf{n}}_m)$ models microphone directivity.

4.3 SPECTRAL SUPERVISION

A key contribution of our framework lies in how we supervise the learning of complex acoustic responses in the frequency domain. Rather than predicting time-domain impulse responses and deriving frequency behavior indirectly—as in prior works—we operate entirely in the spectral domain and define a suite of loss functions that target perceptual alignment, hardware-aware weighting, and physically consistent attenuation. Let $H(f) \in \mathbb{C}$ and $\hat{H}(f) \in \mathbb{C}$ denote the ground-truth and predicted complex frequency responses at a receiver, and let $w(f) \geq 0$ denote a frequency-specific weight that can encode hardware or perceptual importance.

Weighted complex, magnitude, and phase losses. We decompose the spectral supervision into three complementary terms: one for the real and imaginary parts (denoted by Re[.] and Im[.]), one for the magnitude, and one for phase. These jointly ensure accurate complex-valued reconstruction while allowing flexible emphasis through $W_{spec}(f)$, $W_{mag}(f)$ and $W_{phase}(f)$:

$$L_{\text{spec}} = \sum_{f} W_{spec}(f) \left(|\text{Re}[H(f)] - \text{Re}[\hat{H}(f)]| + |\text{Im}[H(f)] - \text{Im}[\hat{H}(f)]| \right), \tag{8}$$

$$L_{\text{mag}} = \sum_{f} W_{mag}(f) \left| |H(f)| - |\hat{H}(f)| \right|, \tag{9}$$

$$L_{\text{phase}} = \sum_{f} W_{phase}(f) \left(|\cos \angle H(f) - \cos \angle \hat{H}(f)| + |\sin \angle H(f) - \sin \angle \hat{H}(f)| \right). \tag{10}$$

These losses provide fine-grained frequency-level control. For example, frequencies in crossover regions of a speaker can be downweighted to avoid unstable learning, while perceptually important midbands can be emphasized.

Spectral envelope smoothing. Acoustic spectra in real-world environments often exhibit narrow-band fluctuations due to interference, which are perceptually less important than the broader spectral shape. Inspired by standard practices in audio engineering and car-cabin equalization, we regularize the predicted and ground-truth log-magnitude spectra while calculating envelope loss L_{env} using an exponential smoothing filter \mathcal{S} :

$$L_{\text{env}} = \sum_{f} \left| \mathcal{S}\left(\log(|H(f)| \cdot w(f) + \epsilon)\right) - \mathcal{S}\left(\log(|\hat{H}(f)| \cdot w(f) + \epsilon)\right) \right|,\tag{11}$$

where ϵ is a small positive constant to avoid singularities. This regularization encourages fidelity to the broadband spectral envelope while tolerating harmless fine-grained ripples, leading to smoother convergence and perceptually cleaner reconstructions.

Kramers–Kronig consistency regularizer. As introduced in Sec. 3.4, in physical media, frequency-dependent attenuation and dispersion are coupled by the Kramers–Kronig (KK) relations. To enforce this constraint in learning, we define:

$$\hat{\beta}(f) = \mathcal{H}\{\sigma\}(f), \qquad L_{KK} = \sum_{f \in \mathcal{B}} \left(\beta(f) - \kappa \,\hat{\beta}(f)\right)^2, \tag{12}$$

where \mathcal{H} is a discrete Hilbert transform implemented via two-sided even extension and a raised-cosine taper, κ is a learnable scalar for scaling alignment, and \mathcal{B} is a frequency band mask to exclude unreliable bins (e.g., DC/Nyquist). This term ensures the learned attenuation field respects causality and avoids unphysical phase artifacts.

Total objective. The complete spectral loss is a weighted combination of the above components:

$$L_{\text{total}} = \lambda_{\text{spec}} L_{\text{spec}} + \lambda_{\text{mag}} L_{\text{mag}} + \lambda_{\text{phase}} L_{\text{phase}} + \lambda_{\text{env}} L_{\text{env}} + \lambda_{\text{KK}} L_{\text{KK}} + L_{\text{aux}}, \tag{13}$$

where $\lambda_{\{\cdot\}}$ are hyperparameters controlling the contribution of each term. L_{aux} denotes auxiliary loss terms carried over from prior work Yamamoto et al. (2020); Majumder et al. (2022) (e.g., multi-resolution STFT loss and energy-shape penalties) and are used for stability rather than driving the primary supervision.

Together, these losses constitute a principled and physically grounded spectral supervision strategy. They allow our model to align with both perceptual and physical constraints—capturing sharp resonances, respecting causal propagation, and adapting to hardware-specific frequency responses—while operating entirely in the frequency domain.

5 EXPERIMENTS

We evaluate INFER on the task of reconstructing the 3D frequency response field inside car cabins. Given measured impulse responses at 48 kHz, the model is trained to predict the complex frequency-domain response at unseen receiver positions. We compare INFER to prior state-of-the-art methods—NAF, INRAS, and AVR—focusing on reconstruction accuracy across frequency bands.

5.1 Datasets

Simulated. We evaluate our method on both simulated and real-world datasets. The simulated data is generated using COMSOL's *Car Cabin Acoustics—Transient Analysis* module, which solves the time-dependent wave equation with realistic, frequency-dependent boundary admittances. We extract impulse responses at 216 receiver positions across the cabin geometry.

Real. For real-world evaluation, we collect measurements in both the BUCK vehicle mock-up and a Tesla Model X using five loudspeakers and a 16-channel UMA-16 microphone array. We record 4096-sample IRs at 48 kHz with physically measured speaker and microphone positions. Fig. 3 shows the data collection environment and hardware for both BUCK(left) and Tesla model X(right).

5.2 IMPLEMENTATION DETAILS

The input to our model consists of a 3D query point $\mathbf{p} \in \mathbb{R}^3$, transmitter location $\mathbf{p}_{tx} \in \mathbb{R}^3$, and directions $(\hat{\mathbf{n}}_{tx}, \hat{\mathbf{n}}) \in \mathbb{R}^3$ representing the emitter and receiver orientations. All input coordinates are







Figure 3: **Data collection setup.** (a)Left: Data is collected in controlled environment - The BUCK, which is a vehicle mockup with realistic car interior and acoustic frontend. (b)Right: Data is also collected in real environment - Tesla Model X.

Table 1: Metrics for BUCK and Tesla setups (lower is better).

		Buck							Tesla							
Method	Amp	Ang	Spec	STFT	Ene.	Env.	T60	EDT	Amp	Ang	Spec	STFT	Ene.	Env.	T60	EDT
INRAS	0.29	1.54	0.9	1.6	7.13	2.79	3.0	3.0	0.43	1.63	1.0	2.2	3.93	3.82	14.6	109.8
NAF	0.14	0.54	0.3	1.0	5.55	1.13	1.3	1.7	0.48	1.63	1.2	2.2	2.25	4.13	10.0	8.1
AVR	0.21	0.81	0.5	1.5	7.95	2.06	3.2	2.4	0.28	1.61	1.0	2.7	5.28	6.92	49.6	24.0
Ours	0.12	0.50	0.2	1.2	5.56	0.95	9.8	2.6	0.14	0.59	0.3	1.0	1.57	1.45	8.4	4.0

encoded using a hash grid based encoding. The model outputs the corresponding density $\sigma[f] \in \mathbb{C}^{\mathbb{T}}$ and frequency signal $H[f] \in \mathbb{C}^{\mathbb{T}}$ at that query point. The model is implemented as MLPs with 6 fully connected layers and 256 hidden units per layer, using ReLU activations. Rendering is performed using ray marching with 64 points per ray, accumulating complex-valued attenuation and delay across the path as described in Sec. 4.2. We integrate over 64×32 azimuth–elevation rays per receiver to form the output signal. We train all models using the Adam optimizer with an initial learning rate of 5×10^{-4} . Training takes approximately 24 hours on a single NVIDIA L40 GPU. All baseline models are trained with the same network size and data splits for fair comparison.

5.3 QUANTITATIVE RESULTS

Across both BUCK and Tesla, INFER achieves the lowest errors on the core frequency-domain metrics (amplitude, angle, spectral, envelope), outperforming NAF, INRAS, and AVR (Table 1). On BUCK, our model is competitive on energy but lags on time-domain reverberation metrics (T60, EDT), which is consistent with our frequency-centric supervision. On Tesla, INFER also leads on the time-domain metrics (lowest T60 and EDT), indicating better generalization. Per-frequency analysis (Table 2) shows consistent gains: INFER attains the best magnitude and phase errors at every reported band, with particularly large phase advantages at low frequencies (e.g., 180 Hz: 0.029 vs. 0.076 for the next best), while maintaining the best average magnitude error. The details on metric calculation can be found in Appendix 8.

5.4 QUALITATIVE RESULTS

Baseline comparison (**Fig. 1**). At 720 Hz the cabin exhibits mid–high modal density with pronounced interference and wrapped phase discontinuities. The left panel (magnitude) shows two dominant high–energy lobes separated by low–energy troughs in the ground truth. Our method reproduces both the location and contrast of these structures, preserving the nodal troughs and avoiding spurious speckle. AVR and NAF show over–textured patterns with high–frequency artifacts and INRAS collapses toward a nearly uniform field. The right panel (phase) smoothly varying phase across two regions. Our prediction aligns with both in terms of the global gradient and the placement, yielding continuous phase evolution. AVR and NAF introduce local phase jitter and INRAS is nearly constant, indicating it does not recover the true phase evolution at this frequency.

Across-frequency behavior (Figs. 4). We evaluate a wide range of frequencies that transition from low-modal regimes to highly complex interference patterns. At 180 Hz the ground truth is

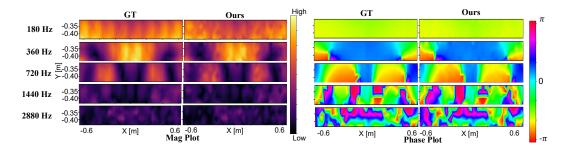


Figure 4: **Qualitative results.** (Left) Spatial plots comparing ground truth (GT), and our method's ability to reconstruct magnitude and phase field across various frequency bands.

Table 2: Per-frequency averages (lower is better). Frequencies are subcolumns of Mag Err and Phase Err: rows are methods.

	Mag Err							Phase Err						
Method	180	360	720	1440	2880	Avg	180	360	720	1440	2880	Avg		
INRAS	0.751	0.602	0.662	0.384	0.280	0.536	0.100	0.824	1.232	1.285	1.376	0.963		
NAF	0.331	0.277	0.242	0.154	0.150	0.231	0.076	0.284	0.398	0.500	0.413	0.334		
AVR	0.465	0.388	0.255	0.183	0.217	0.302	0.105	0.187	0.151	0.330	0.665	0.288		
Ours	0.149	0.152	0.125	0.154	0.118	0.140	0.029	0.076	0.081	0.194	0.322	0.140		

dominated by a slowly varying field; our model reproduces this near-uniform phase and the smooth, weakly varying magnitude without introducing spurious structure. At 360–720 Hz, where the cabin begins to exhibit distinct standing-wave patterns, our reconstructions recover both the placement and contrast of high/low-energy lobes together with the associated phase gradients and wrap seams. At higher frequencies (1440–2880 Hz) the field contains rapid spatial oscillations and multiple discontinuities. Despite this, our predictions remain stable: magnitude maps retain fine-scale contrast without speckle or over-smoothing, and phase maps capture the correct number and placement of wraps with coherent local gradients.

6 DISCUSSION AND FUTURE WORK

This work introduces INFER, a novel spectral-domain neural representation for car cabin acoustics, enabling accurate frequency response reconstruction from sparse measurements. Our method surpasses prior baselines in both magnitude and phase fidelity, and remains physically grounded through causality-aware regularization. Beyond its immediate impact on spatial audio modeling and personalization, our approach opens avenues for integrating learned acoustic fields into downstream tasks such as adaptive ANC, directional speech enhancement, and real-time audio rendering. Future extensions include joint modeling across passenger positions, integrating speaker-specific transfer functions, and exploring generalization to unseen vehicle geometries or dynamic cabin conditions.

7 CONCLUSION

We introduced INFER, a novel spectral-domain framework that models acoustic propagation in confined environments using implicit neural representations. By operating directly in the frequency domain, our method enables perceptually grounded supervision, hardware-aware weighting, and physically consistent regularization through the Kramers–Kronig constraint. Our differentiable renderer explicitly accounts for phase and attenuation via complex-valued ray integration, yielding spatially coherent and frequency-resolved reconstructions. Extensive evaluations on real and simulated car cabin datasets demonstrate that INFER substantially outperforms prior time-domain and hybrid approaches, achieving over 50% improvement in phase accuracy and 39% in magnitude fidelity relative to the best baseline. We believe this work paves the way for accurate and physically grounded neural acoustic modeling in automotive spaces.

8 REPRODUCIBILITY

Anonymized code and demo datasets will be available on our webpage (https://anonymous1415510-spec.github.io). We provide details about comparison with other algorithms to facilitate reproducing our results. All details about the hyperparameters, environment specifications, and real-world experiment setup are provided in the appendix or the website.

REFERENCES

- Giampiero Accardo, Paolo Chiariotti, Bram Cornelis, Mahmoud El-Kafafy, B Peeters, K Janssens, and Milena Martarelli. Experimental acoustic modal analysis of an automotive cabin: challenges and solutions. In *Journal of Physics: Conference Series*, volume 1075, pp. 012026. IOP Publishing, 2018.
- John T Etgen and Michael J O'Brien. Computational methods for large-scale 3d acoustic finite-difference modeling: A tutorial. *Geophysics*, 72(5):SM223–SM230, 2007.
- Angelo Farina, Emanuele Ugolotti, et al. Spatial equalization of sound systems in cars. In *Proc. of* 15th AES Conference" Audio, Acoustics & Small Spaces", Copenhagen, Denmark, volume 31, pp. 10–2, 1998.
- Jinlong Feng, Xiaoping Zheng, Haitao Wang, HongTao Wang, Yuanjie Zou, Yinghua Liu, and Zhenhan Yao. Low-frequency acoustic-structure analysis using coupled fem-bem method. *Mathematical Problems in Engineering*, 2013(1):583079, 2013.
- Philipp Holl, Vladlen Koltun, and Nils Thuerey. Learning to control pdes with differentiable physics. *arXiv* preprint arXiv:2001.07457, 2020.
- Shoichi Koyama, Enzo De Sena, Prasanga Samarasinghe, Mark RP Thomas, and Fabio Antonacci. Past, present, and future of spatial audio and room acoustics. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2025.
- Heinrich Kuttruff. Room acoustics. Crc Press, 2016.
- Zitong Lan, Chenhao Zheng, Zhiwei Zheng, and Mingmin Zhao. Acoustic volume rendering for neural impulse response fields. *Advances in Neural Information Processing Systems*, 37:44600–44623, 2024.
- Ren-Zun Lian. Generalized huygens–fresnel principle. In 2023 9th International Conference on Computer and Communications (ICCC), pp. 1616–1627. IEEE, 2023.
- Susan Liang, Chao Huang, Yapeng Tian, Anurag Kumar, and Chenliang Xu. Av-nerf: Learning neural fields for real-world audio-visual scene synthesis. *Advances in Neural Information Processing Systems*, 36:37472–37490, 2023.
- Yanshan Liu, Xiangyang Zeng, and Haitao Wang. 3d meshless fem-bem model for prediction of sound fields in cabins due to external sound disturbances. *Journal of Vibroengineering*, 19(7): 5346–5360, 2017.
- Andrew Luo, Yilun Du, Michael Tarr, Josh Tenenbaum, Antonio Torralba, and Chuang Gan. Learning neural acoustic fields. *Advances in Neural Information Processing Systems*, 35:3165–3177, 2022.
- Sagnik Majumder, Changan Chen, Ziad Al-Halah, and Kristen Grauman. Few-shot audio-visual learning of environment acoustics. *Advances in Neural Information Processing Systems*, 35: 2522–2536, 2022.
- Julien NP Martel, David B Lindell, Connor Z Lin, Eric R Chan, Marco Monteiro, and Gordon Wetzstein. Acorn: Adaptive coordinate networks for neural scene representation. *arXiv* preprint *arXiv*:2105.02788, 2021.

543

544

546 547

548

549

550 551

552

553

554

555

556

558

559

560

561

562 563

564

565

566

567 568

569

570

571

572

573

574

575

576

577

578 579

580

581

582

583

584

585 586

587

588

590 591

592

593

- 540 Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. Communications 542 of the ACM, 65(1):99–106, 2021.
 - Qi Mo, Hengchin Yeh, and Dinesh Manocha. Tracing analytic ray curves for light and sound propagation in non-linear media. IEEE transactions on visualization and computer graphics, 22(11): 2493-2506, 2015.
 - Amirali Molaei, Amirhossein Aminimehr, Armin Tavakoli, Amirhossein Kazerouni, Bobby Azad, Reza Azad, and Dorit Merhof. Implicit neural representation in medical imaging: A comparative survey. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 2381– 2391, 2023.
 - Matthew ODONNELL, ET Jaynes, and JG Miller. Kramers-kronig relationship between ultrasonic attenuation and phase velocity. Acoustical Society of America, Journal, 69:696–701, 1981.
 - Tribhuvanesh Orekondy, Pratik Kumar, Shreya Kadambi, Hao Ye, Joseph Soriaga, and Arash Behboodi. Winert: Towards neural ray tracing for wireless channel modelling and differentiable simulations. In The Eleventh International Conference on Learning Representations, 2023.
 - Andrew J Oxenham. How we hear: The perception and neural coding of sound. Annual review of psychology, 69(1):27-50, 2018.
 - Vincent Sitzmann, Julien N. P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. ArXiv, abs/2006.09661, 2020. URL https://api.semanticscholar.org/CorpusID:219720931.
 - Kun Su, Mingfei Chen, and Eli Shlizerman. Inras: Implicit neural representation for audio scenes. Advances in Neural Information Processing Systems, 35:8144–8158, 2022.
 - Harshvardhan Takawale and Nirupam Roy. Spinrv2: Implicit neural representation for passband fmcw radars, 2025.
 - Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. Advances in neural information processing systems, 33:7537-7547, 2020.
 - Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. arXiv preprint arXiv:2106.10689, 2021.
 - YiPing Wang, Xin Zhen, Jing Wu, ZhengQi Gu, ZhenYi Xiao, and Xue Yang. Hybrid cfd/fembem simulation of cabin aerodynamic noise for vehicles traveling at high speed. Science China *Technological Sciences*, 56(7):1697–1708, 2013.
 - Ryuichi Yamamoto, Eunwoo Song, and Jae-Min Kim. Parallel wavegan: A fast waveform generation model based on generative adversarial networks with multi-resolution spectrogram. In ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 6199–6203. IEEE, 2020.
 - T Yoshimura, M Saito, S Maruyama, and S Iba. Modal analysis of automotive cabin by multiple acoustic excitation. In *Proceedings of ISMA*, volume 6880, 2012.
 - Xinyun Zhang, Ruiqi Yu, and Shuang Ren. Neural implicit representations for multi-view surface reconstruction: A survey. IEEE Transactions on Visualization and Computer Graphics, 2025.

DISCLOSURE OF LLM USAGE FOR WRITING

Large Language Models (LLMs) were used solely for grammar refinement and polishing of the manuscript text. All ideas, technical content, experimental design, and analysis were independently developed by the authors without LLM assistance.

TRAINING DETAILS AND REPRODUCIBILITY

A.1 REPRODUCIBILITY

Anonymized code and demo datasets will be available on our webpage (https://anonymous1415510-spec.github.io). We provide details about comparison with other algorithms to facilitate reproducing our results. All details about the hyperparameters, environment specifications, and real-world experiment setup are provided in the appendix or the website.

A.2 MODEL ARCHITECTURE AND RENDERER DETAILS

We adopt the AVRModel_complex_FD_FreqDep_PhaseCorrection model architecture with the AVRRenderFD_FreqDep_PhaseCorrection_KK renderer. This setup is designed for complex-valued frequency domain rendering and enables physically grounded learning with explicit modeling of attenuation and phase velocity.

Key architectural components:

- Complex Frequency Field Prediction: The model learns frequency-dependent attenuation fields $\delta(f) = \sigma(f) + j\beta(f)$ and complex responses H(f) using MLPs operating on hash-encoded spatial coordinates.
- **Separate Signal and Attenuation Networks**: Frequency-specific signal and attenuation values are predicted using distinct encoders and MLPs.
- Directional Encodings: Transmitter and receiver directions are encoded using spherical harmonics.
- Renderer Pipeline: Rays are sampled in spherical directions with integration over 64 azimuth × 32 elevation rays, each with 64 samples from near=0 to far=4 meters. Cumulative attenuation is applied using $\exp(-\sum \sigma_i \Delta u_i + j \sum \beta_i \Delta u_i)$.

Parameter	Value							
Learning Rate	5×10^{-4} (cosine annealing to 5×10^{-5})							
Optimizer	Adam							
Total Iterations	15,000							
Batch Size	1							
Rendering Samples	64 per ray							
Azimuth × Elevation Rays	64×32							
Speed of Sound	343.8 m/s							
Sampling Frequency	48,000 Hz							
Path Loss Exponent	1							
Layers	8 fully connected layers							
Hidden Units	256 neurons per layer							
Activation	ReLU							
Positional Encoding	10 frequencies for spatial and directional input							

Table 3: Training hyperparameters and network architecture for INFER.

A.3 Loss Functions and Weights

The total training objective is composed of multiple terms designed to supervise the model's output across spectral amplitude, phase, energy structure, and physical consistency. The overall loss is expressed as:

$$L_{\rm total} = \lambda_{\rm spec} L_{\rm spec} + \lambda_{\rm mag} L_{\rm mag} + \lambda_{\rm phase} L_{\rm phase} + \lambda_{\rm env} L_{\rm env} + \lambda_{\rm energy} L_{\rm energy} + \lambda_{\rm KK} L_{\rm KK} + \lambda_{\rm STFT} L_{\rm MR-STFT}. \tag{14}$$

In our experiments, we set $\lambda_{\text{spec}} = 16$, $\lambda_{\text{mag}} = 4$, $\lambda_{\text{phase}} = 1$, $\lambda_{\text{env}} = 0.25$, $\lambda_{\text{energy}} = 2$, $\lambda_{\text{KK}} = 0.25$, and $\lambda_{\text{STFT}} = 0.25$. Beyond these global weights, we apply frequency-dependent weighting: for

 $L_{
m phase}$ we emphasize low and mid frequencies by setting w(f)=1.2 up to 1.5 kHz (125 bins), w(f)=1 until 5 kHz (425 bins), and then smoothly tapering to 0.8 across the log-frequency axis; for $L_{
m mag}$, weights are 1 up to 1.5 kHz, 1.25 until 5 kHz, and then tapered to 0.8; and for $L_{
m spec}$, we assign 1.25 up to 5 kHz before tapering. These schedules follow psychoacoustic informed weighting strategies, prioritizing perceptually critical bands while deemphasizing unreliable high-frequency bins.

A.4 TRAINING PIPELINE

 We use the script avr_runner_complex_FD_kk.py for training. Each step involves:

- 1. Ray-based spherical integration using normalized receiver and transmitter coordinates.
- 2. Prediction of complex signals and attenuation fields.
- 3. Loss computation including frequency-weighted spectrum and KK regularization.
- 4. Gradient backpropagation with NaN filtering and norm clipping.
- 5. Mixed precision training and GPU memory optimization.

Special Considerations:

- Gradient clipping to max norm 1.
- Automatic mixed precision (AMP) to save memory.
- Complex loss handling via separate $\Re[H(f)]$ and $\Im[H(f)]$ paths.
- KK regularizer computed using discrete Hilbert transform with frequency masking and tapering.

A.5 REPRODUCIBILITY CHECKLIST

Software Environment:

- Python 3.8, PyTorch 1.12.0, CUDA 11.6
- tinycudann 1.6, auraloss 0.4.0, tensorboard 2.8.0

Training Script (Single GPU):

```
python avr_runner_complex_FD_kk.py \
    --config config_files/avr_buck_complex_dir_FD.yml \
    --model_type AVRModel_complex_FD_FreqDep_PhaseCorrection \
    --renderer_type AVRRenderFD_FreqDep_PhaseCorrection_KK \
    --batchsize 1 \
    --dataset_dir /path/to/dataset
```

A.6 HARDWARE REQUIREMENTS

Minimum:

• GPU: NVIDIA RTXA6000 (10GB+ VRAM)

Recommended:

• GPU: L40S

Training time is approximately 24 hours on a single L40S.

A.7 EVALUATION METRICS

- Spectral Distance (SD): $20 \log_{10}(|||H(f)| |\hat{H}(f)|||/|||\hat{H}(f)|||)$
- Phase Error (PRE): $\|\angle H(f) \angle \hat{H}(f)\|_2/\pi$
- Energy Ratio, Cumulative Energy Error in frequency domain
- Time-Domain Metrics: PSNR, MR-STFT loss
- KK Violation Metric: $\|\beta(f) \mathcal{H}[\sigma(f)]\|/\|\beta(f)\|$

A.8 AUDIO HARDWARE SPECIFICATION

We detail here the acoustic transducer setup used for our data collection in the *BUCK* testbed and the production *Tesla Model X* vehicle. Both systems were equipped with a rich spatial arrangement of loudspeakers and a high-fidelity microphone array to facilitate spatial audio capture and reconstruction.

Speaker Configuration. While the Tesla Model X uses the default speakers, In BUCK, the active speakers used for sound excitation include:

- Center Dash Speaker: A 3.5-inch wideband driver, such as the SLA Ram3 or Dayton Audio DMA90-4, capable of full-range output from 85.00 Hz to 20.00 kHz. In typical configurations, these are high-passed at approximately 100.00 Hz to avoid low-frequency distortion.
- Rear Door Speakers: Morel Tempo Ultra Integra 402 or 602 coaxial hybrids with wideband support (55.00 Hz to 22.00 kHz), high sensitivity, and power handling up to 120.00 WRMS. These speakers internally crossover between woofer and tweeter around 2.50 kHz–3.00 kHz.
- **Rear Height Speakers:** Tang Band T2-2136SF full-range modules and Morel CCWR254 midrange drivers, mounted in ceiling/rear hatch positions to introduce vertical spatial content, spanning 80.00 Hz to 20.00 kHz. Crossover filters are typically applied around 800.00 Hz-1.00 kHz depending on system design and companion driver.

This layout approximates a 7.1.4 immersive audio setup and enables extensive sampling of reverberant and directional field responses across both testbeds.

Microphone Array. We use the commercially available **MiniDSP UMA-16** USB microphone array, which offers 16 omnidirectional MEMS microphones in a linear array form factor. This array supports high-resolution spatial sampling across the cabin, enabling dense reconstruction of directional impulse responses.

A.9 BASELINE METHODS

To rigorously evaluate the effectiveness of our proposed system INFER, we compare against three representative baselines, each reflecting a different class of acoustic modeling approach:

- AVR: A hybrid time–frequency domain neural field that learns time-domain impulse responses via a differentiable renderer. While AVR applies frequency-domain path delays in its rendering, the model is supervised in the time domain and does not explicitly learn frequency-dependent attenuation or dispersion.
- NAF (Neural Acoustic Field): A neural field trained directly in the time domain using MSE and time-domain perceptual losses. NAF ignores frequency-domain supervision and is evaluated primarily on time-domain waveform fidelity.
- INRAS (Impulse Response as Signal): A signal regression approach where the model directly regresses to the complex impulse response waveform as a 1D signal. INRAS uses STFT-based perceptual loss, but it does not exploit any spatial priors or directional conditioning.

Each baseline is re-implemented in our codebase with their respective loss functions and evaluation metrics faithfully reproduced, using the same training datasets, preprocessing pipelines, and neural architecture backbones where applicable.

A.10 Training Configuration

All baselines are trained with identical configurations to ensure fair comparisons:

- **Dataset:** We use the same training/validation/test splits from our real (BUCK) and synthetic (Tesla) datasets for all methods.
- **Resolution:** The frequency bins, spatial sampling resolution, and directional integration are matched across all methods.
- Training Epochs: All models are trained for 500 epochs with early stopping based on validation loss.
- Batch Size: Batch size of 1 is used due to GPU memory constraints, consistent with prior volumetric rendering works.
- Evaluation: All comparisons are evaluated on both magnitude and phase accuracy across the frequency range of interest, in addition to perceptual STFT loss and energy-based metrics.

A.11 Loss Function Implementation

AVR. We follow the original AVR formulation and use the same set of losses described in the paper.

NAF. The NAF baseline is trained using the standard losses introduced in its original work, without any additional frequency-domain regularization.

INRAS. For INRAS, we adopt the exact losses specified in the original paper, without modification.

A.12 ARCHITECTURAL MODIFICATIONS

To isolate the effects of spectral supervision and renderer formulation, all baseline models are built upon the same backbone MLP architecture as our method:

- 8-layer fully connected network with sinusoidal positional encoding.
- Input: $(\mathbf{p}_{tx}, \hat{\mathbf{n}}_{tx}, \mathbf{x}, \hat{\mathbf{n}})$ with appropriate frequency and spatial encodings.
- Output: Real-valued waveform or complex spectrum depending on method.

A.13 EVALUATION METRICS

We evaluate all baselines and our method using the following metrics:

- **Amplitude Error**: L1 distance between predicted and ground-truth magnitude spectra across frequency.
- Phase Error: Cosine-sine angular difference between predicted and true phase.
- Envelope Error: Smoothed log-magnitude deviation (via exponential smoothing).
- MR-STFT Loss: Multi-resolution STFT loss to evaluate perceptual reconstruction in time domain.
- **Energy Accuracy**: Deviation from cumulative spectral energy profiles.

A.14 Notes on Fairness and Robustness

To ensure fairness in evaluation:

- All models are trained with the same GPU hardware, random seed initialization, and Py-Torch version.
- We use the same optimizer (Adam) and learning rate schedule across all models unless otherwise noted.
- All baselines are evaluated using our standardized renderer and metric pipeline to eliminate post-processing inconsistencies.
- We tune loss weights and learning rates for each baseline to ensure their best performance under our training conditions.

Overall, our comparison demonstrates that INFER substantially outperforms these baselines across spectral and perceptual metrics due to its physics-informed supervision and frequency-aware modeling.