FuseCodec: Semantic-Contextual Fusion and Supervision for Neural Codecs

Anonymous ACL submission

Abstract

Speech tokenization enables discrete representation and facilitates speech language modeling. However, existing neural codecs capture lowlevel acoustic features, overlooking the 006 semantic and contextual cues inherent to human speech. While recent efforts introduced 007 semantic representations from self-supervised speech models or incorporated contextual representations from pre-trained language models, 011 challenges remain in aligning and unifying the semantic and contextual representations. We in-012 troduce FuseCodec, which unifies acoustic, semantic, and contextual representations through strong cross-modal alignment and globally informed supervision. We propose three complementary techniques: (i) Latent Representation 017 Fusion, integrating semantic and contextual fea-019 tures directly into the encoder latent space for robust and unified representation learning; (ii) Global Semantic-Contextual Supervision, supervising discrete tokens with globally pooled 023 and broadcasted representations to enhance temporal consistency and cross-modal alignment; and (iii) Temporally Aligned Contextual Supervision, strengthening alignment by 027 dynamically matching contextual and speech tokens within a local window for fine-grained token-level supervision. We further introduce FuseCodec-TTS, demonstrating our methodology's applicability to zero-shot speech synthesis. Empirically, FuseCodec achieves state-of-the-art performance in LibriSpeech, surpassing EnCodec, SpeechTokenizer, and DAC in transcription accuracy, perceptual quality, intelligibility, and speaker similarity. 037 Results highlight the effectiveness of contextually and semantically guided tokenization for speech tokenization and downstream tasks.

1 Introduction

043

Tokenization has become foundational in natural language processing (NLP), enabling language models to learn discrete representations, while facilitating efficient autoregressive modeling and scalable downstream applications (Schmidt et al., 2024). Inspired by this paradigm, the speech domain has increasingly adopted neural codecs, popularized by Encodec (Défossez et al., 2022) and SoundStream (Zeghidour et al., 2022). Neural codecs tokenize speech using an encoder, residual vector quantizer, and decoder architecture, enabling modeling discrete representations suitable for modular extension to downstream tasks such as speech synthesis (Wang et al., 2023). 044

045

046

047

051

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

081

However, the continuous and multidimensional nature of human speech makes learning discrete representations inherently challenging (Ju et al., 2024). While neural codecs learn acoustic representations (waveform and low-level signal characteristics), they struggle to capture high-level semantics requiring downstream models to adopt additional self-supervised masked language objectives to derive semantic representations (phonetic content and linguistic meaning) (Borsos et al., 2023). To address this drawback, recent neural codec architectures incorporated semantic distillation from pretrained self-supervised speech models (Zhang et al., 2024; Défossez et al., 2024), improving the quality of speech reconstruction and the semantic aspect of learned representations.

In addition, another fundamental aspect of human speech remains missing in above mentioned works: speech is inherently grounded in context and surrounding cues (Brown et al., 2022). Discrete speech representations, lacking grounding in context, fall short of capturing this essential attribute (Hallap et al., 2023). While language models have demonstrated strong capabilities in learning such contextual dependencies from text corpora (Devlin et al., 2019a; Peters et al., 2018), speech tokenizers have yet to fully leverage these capabilities. Although a recent neural codec (Ahasan et al., 2024) explored matching discrete speech representations with contextual representation from a pre-trained

125 126

127

128 129

130

132

133

134

135

136

language model, it falls short in effective crossmodal alignment, constraining the model's ability to fully unify semantic and contextual information.

Therefore, despite recent advances, several challenges remain unaddressed. Firstly, current approaches fail to unify all three aspects of discrete speech representation: acoustic (learned by neural codecs), semantic (from self-supervised speech models), and contextual (from language models). Most work incorporates only semantic information (Zhang et al., 2024; Défossez et al., 2024; Ye et al., 2024), neglecting contextual grounding. Secondly, while a recent effort (Ahasan et al., 2024) attempts to integrate contextual representations, it lacks effective mechanisms for aligning text and speech modalities. Thirdly, existing methods rely on similarity-based objectives for representation matching without directly incorporating information into the latent space, limiting coherence and downstream performance. We address these challenges through our proposed methodologies, while preserving the core architecture and utilizing frozen representations with zero inference overhead.

To address these challenges, we propose a speech tokenization framework with three different strategies/variations that enrich discrete speech representations with unified and aligned semantic and contextual information. Our first strategy involves (i) Latent Representation Fusion, which integrates semantic and contextual embeddings into the encoder's latent space through cross-modal attention and additive fusion, resulting in more robust and coherent representations. Building on this, we present (ii) Global Semantic-Contextual Supervision, where globally pooled and broadcasted modality vectors supervise each quantized token across time, facilitating temporally consistent and globally informed representation learning. To enforce explicit alignment, we introduce another strategy: (iii) Temporally Aligned Contextual Supervision, which dynamically matches contextual and speech tokens prior to timestep-level similarity supervision, enabling fine-grained cross-modal alignment and enhancing representation quality.

Then, we instantiate our framework through three model variants: FuseCodec-Fusion with Latent Representation Fusion, FuseCodec-Distill with Global Semantic-Contextual Supervision, and FuseCodec-ContextAlign with Temporally Aligned Contextual Supervision. FuseCodec establishes new state-of-the-art performance on the LibriSpeech test set (Panayotov et al., 2015) by integrating contextual and semantic guidance into the learning of discrete speech tokens. Specifically, FuseCodec-Fusion achieves the best scores 139 in transcription accuracy (WER 3.99, WIL 6.45), 140 intelligibility (STOI 0.95), and perceptual quality 141 (ViSQOL 3.47, PESQ 3.13), outperforming En-142 Codec (Défossez et al., 2022), SpeechTokenizer 143 (Zhang et al., 2024), and DM-Codec (Ahasan 144 et al., 2024). FuseCodec-Distill further achieves 145 the highest UTMOS (3.65) and speaker similar-146 ity (0.996), highlighting its strength in percep-147 tual naturalness and speaker fidelity. Meanwhile, 148 FuseCodec-ContextAlign provides a strong trade-149 off between interpretability and performance, with 150 particularly competitive scores in UTMOS (3.65) 151 and similarity (0.995). These results underscore 152 the effectiveness of incorporating contextual and 153 semantic signals into the tokenization process for 154 high-quality speech reconstruction. 155

137

138

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

Therefore, our key contributions are:

- We introduce three novel neural codecs based on our method: Latent Representation Fusion (FuseCodec-Fusion), Global Semantic-Contextual Supervision (FuseCodec-Distill), and Temporally Aligned Contextual Supervision (FuseCodec-ContextAlign).
- Our framework tackles different limitations of neural codecs by integrating semantic and contextual information through distinct methods, improving cross-modal alignment and enhancing discrete representation learning.
- We demonstrate the utility of our approach in a downstream TTS model and validate each component with extensive ablation studies.
- FuseCodec achieves state-of-the-art performance on LibriSpeech reducing transcription error and improving speech naturalness.

2 Proposed Method

As shown in Figure 1, we first introduce the speech discretization pipeline (§2.1) and describe the extraction of semantic and contextual representations from pre-trained models (§2.2). We then present three strategies for integrating multimodal guidance into speech tokenization: (i) Latent Representation Fusion (§2.3.1), (ii) Global Semantic-Contextual Supervision (§2.3.2), and (iii) Temporally Aligned Contextual Supervision (§2.3.3). Finally, we outline the training objective $(\S2.4)$ and the extension to a text-to-speech task ($\S2.5$).



Figure 1: Overview of the FuseCodec speech tokenization framework. Input speech x is encoded into latent features \mathbf{Z} , then quantized into discrete tokens $\mathbf{Q}^{(1:K)}$ via residual vector quantization (RVQ). To enrich these tokens, we incorporate semantic $(\mathbf{S}_i, \hat{\mathbf{S}})$ and contextual $(\mathbf{C}_i, \hat{\mathbf{C}}, \mathbf{C}^*)$ representations from frozen pre-trained models. Global vectors $\hat{\mathbf{S}}$ and $\hat{\mathbf{C}}$ are formed via mean pooling and [CLS] selection, respectively. We propose three strategies: (i) Latent Representation Fusion, injecting global vectors $\hat{\mathbf{S}}, \hat{\mathbf{C}}$ with Z to yield fused latent \mathbf{Z}' ; (ii) Global Semantic-Contextual Supervision, supervising $\mathbf{Q}^{(1)}$ with global vectors; and (iii) Temporally Aligned Contextual Supervision, aligning full contextual embeddings $\{\mathbf{C}_i\}$ to RVQ outputs via a windowed matching algorithm to form \mathbf{C}^* .

2.1 Discrete Speech Representation

186

187

188

194

195

196

197

199

200

203

207

Discrete tokens serve as the foundation of neural codec-based speech-language models. Following established approaches (Défossez et al., 2022; Zhang et al., 2024; Ahasan et al., 2024), we discretize audio using an encoder-quantizer setup.

Given an input speech waveform x, an encoder *E* compresses x into a sequence of latent representations $\mathbf{Z} = {\{\mathbf{z}_i\}}_{i=1}^{T'}$, where *T'* is the number of encoded frames. The encoder output **Z** is then passed through a Residual Vector Quantization module (RVQ), consisting of *K* quantization layers. Each layer *k* produces a sequence of token indices ${\{q_i^{(k)}\}}_{i=1}^{T'}$. For each token index $q_i^{(k)}$, we retrieve its corresponding embedding from the *k*-th codebook, resulting in a sequence of quantized vectors $\mathbf{Q}^{(k)} = {\{\mathbf{q}_i^{(k)}\}}_{i=1}^{T'}$, where $\mathbf{q}_i^{(k)} \in \mathbb{R}^D$, with *D* denoting the embedding dimensionality. We use the embeddings from the first quantization layer $\mathbf{Q}^{(1)}$ as the discrete representation of speech, guided with multimodal representations.

2.2 Multimodal Representation Extraction

Concurrently, we extract representations from pretrained models. Specifically, we obtain contextual
representations from a pre-trained language model,
which are dynamic, token-level embeddings that
adapt to surrounding text (Devlin et al., 2019b; Peters et al., 2018). In parallel, we derive semantic

representations from a pre-trained self-supervised speech model, which capture the high-level structure and meaning (Borsos et al., 2023).

Contextual Representation. The input speech waveform **x** is transcribed into text **x'** using a pre-trained Automatic Speech Recognition (ASR) model A, such that $\mathbf{x}' = A(\mathbf{x})$. The ASR model functions purely as a speech-to-text converter and remains detached during training. The transcribed text \mathbf{x}' is processed by a pre-trained language model B, which produces a token sequence $\{c_i\}_{i=1}^n$. For each token c_i , we extract hidden states from all L layers, represented as $\{\mathbf{h}_i^{(l)}\}_{l=1}^L$. These are averaged to produce contextual embeddings: $\mathbf{C}_i = \frac{1}{L} \sum_{l=1}^L \mathbf{h}_i^{(l)}$, where $\mathbf{C}_i \in \mathbb{R}^{D'}$, and D' denotes the hidden dimension of the language model.

Semantic Representation. The input speech waveform **x** is passed through a pre-trained self-supervised speech model H, which outputs a sequence of frame-level tokens $\{s_i\}_{i=1}^m$, For each frame s_i , we extract hidden states from all L layers: $\{\mathbf{h}_i^{(l)}\}_{l=1}^L$. These are averaged to obtain semantic embeddings: $\mathbf{S}_i = \frac{1}{L} \sum_{l=1}^L \mathbf{h}_i^{(l)}$, where $\mathbf{S}_i \in \mathbb{R}^{D'}$, and D' denotes the hidden dimension.

2.3 Semantic-Contextual Guidance

Our goal is to enrich discrete speech representations by integrating contextual and semantic information, enabling tighter alignment between 214

215

216

217

218

219

220

221

222

224

225

228

229

232

233

234

235

236

238

239

240

acoustic structure and linguistic meaning. Prior work has explored similar directions: Zhang et al. (2024); Défossez et al. (2024) aligned HuBERTbased semantic features with the first RVQ layer using cosine similarity, while Ahasan et al. (2024) matched BERT-based embeddings to RVQ outputs via padded sequences and similarity loss. However, these methods either rely on a single modality (semantic in Zhang et al. (2024); Défossez et al. (2024)) or lack robust cross-modal alignment (misaligned context in Ahasan et al. (2024)).

242

243

244

246

247

251

256

261

262

263

265

267

269

270

273

274

275

276

277

278

283

285

In contrast, we unify semantic and contextual representations while ensuring robust alignment. For this, we propose three strategies: (i) Latent Representation Fusion (§2.3.1), (ii) Global Semantic-Contextual Supervision (§2.3.2), and (iii) Temporally Aligned Contextual Supervision (§2.3.3)

2.3.1 Latent Representation Fusion

We first propose fusing semantic and contextual representations with the encoder's latent output. The enhanced latents are then passed to the residual vector quantization (RVQ) module, enabling the learning of discrete codes enriched with semantic and contextual information.

We begin by obtaining global semantic and contextual representations. Specifically, we take the average of semantic embeddings $\{\mathbf{S}_i\}_{i=1}^m$ to compute the global semantic vector $\hat{\mathbf{S}} = \frac{1}{m} \sum_{i=1}^m \mathbf{S}_i$. For the textual modality, we select the [CLS] token embedding from the contextual representations $\{\mathbf{C}_i\}_{i=1}^n$, yielding $\hat{\mathbf{C}} = \mathbf{C}_{[CLS]}$.

We then broadcast each global vector across the discrete token sequence length T', forming: $\tilde{\mathbf{S}} = {\{\hat{\mathbf{S}}\}_{t=1}^{T'}, \text{ and } \tilde{\mathbf{C}} = {\{\hat{\mathbf{C}}\}_{t=1}^{T'}. \text{ Broadcasting allows each token to inherit the full semantic or contextual knowledge of the sequence, ensuring every position is enriched with the most informative signal for cross-modal fusion or distillation.$

Next, we apply multi-head cross-attention to enable cross-modal interaction, followed by an MLP projection to match the encoder dimension D:

$$\mathbf{S}' = \operatorname{CrossAttention}(\tilde{\mathbf{S}}, \tilde{\mathbf{C}}, \tilde{\mathbf{C}}) \mathbf{W}_{S}, \\ \mathbf{C}' = \operatorname{CrossAttention}(\tilde{\mathbf{C}}, \tilde{\mathbf{S}}, \tilde{\mathbf{S}}) \mathbf{W}_{C},$$
(1)

where $\mathbf{W}_S, \mathbf{W}_C \in \mathbb{R}^{D' \times D}$ are learned projection matrices and CrossAttention(\cdot) denotes multihead cross-attention. Finally, we fuse the modality signals with the latent representation $\mathbf{Z} \in \mathbb{R}^{T' \times D}$ via additive fusion and modality dropout:

$$\mathbf{Z}' = \mathbf{Z} + (\mathbf{S}' \odot \mathcal{D}_S) + (\mathbf{C}' \odot \mathcal{D}_C), \quad (2)$$

where $\mathcal{D}_S, \mathcal{D}_C \in \{0,1\}^{T' \times D}$ are stochastic dropout masks applied during training. Dropout promotes robustness by preventing the quantized representations from over-relying on the fused modalities (Hussen Abdelaziz et al., 2020), and allows inference using only the encoder signal. The resulting fused representation \mathbf{Z}' is then passed to the RVQ module for discrete speech quantization.

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

330

331

332

333

334

335

2.3.2 Semantic-Contextual Supervision

In addition to latent fusion, we explore an alternative representation supervision strategy, motivated by its effectiveness of similarity matching in prior speech tokenization work (Zhang et al., 2024; Défossez et al., 2024; Ahasan et al., 2024). Unlike previous methods that supervise over feature dimensions or require local frame-level alignment, we introduce a global-to-local time-axis distillation scheme. Specifically, we use global semantic \hat{S} and contextual \hat{C} vectors to supervise the RVQ output across time. This provides temporally consistent guidance and encourages the quantized space to capture modality-aware temporal dynamics.

We adapt the *combined distillation loss* from Ahasan et al. (2024), proposing it to operate along the temporal axis rather than the feature axis. This modification enables more effective alignment of discrete latent representations with temporally distributed semantic and contextual signals, enhancing cross-modal coherence over time.

Given the broadcasted global signals (see 2.3.1) $\tilde{\mathbf{S}}, \tilde{\mathbf{C}} \in \mathbb{R}^{T' \times D'}$, we apply a linear projection to the first-layer RVQ output $\mathbf{Q}^{(1)} \in \mathbb{R}^{T' \times D}$ to align dimensionality: $\mathbf{Q}'^{(1)} = \mathbf{Q}^{(1)}\mathbf{W}$, where $\mathbf{W} \in \mathbb{R}^{D \times D'}$. Finally, we apply semantic-contextual supervision using a temporally-aware distillation loss.

$$\mathcal{L}_{\text{distill}} = -\frac{1}{T'} \sum_{t=1}^{T'} \log \sigma \left(\frac{1}{2} \left[\cos \left(\mathbf{Q}_t^{\prime(1)}, \tilde{\mathbf{S}}_t \right) + \cos \left(\mathbf{Q}_t^{\prime(1)}, \tilde{\mathbf{C}}_t \right) \right] \right),$$
(3)

where $\sigma(\cdot)$ is the sigmoid function and $\cos(\cdot, \cdot)$ denotes cosine similarity. This formulation provides fine-grained temporal supervision using global modality signals, enhancing the representational quality of the learned discrete tokens.

2.3.3 Aligned Contextual Supervision

Building on our use of the global contextual vector $\hat{\mathbf{C}}$ for supervision, we propose a finer-grained approach that leverages the full sequence of contextual embeddings $\{\mathbf{C}_i\}_{i=1}^n$ to supervise the RVQ to-

ken sequence $\{\mathbf{q}_t^{(1)}\}_{t=1}^{T'}$, enabling richer, timesteplevel guidance. A key challenge, however, is the mismatch in sequence lengths between the contextual embeddings (n) and the RVQ output (T').

To address this, we propose a dynamic windowbased alignment strategy that assigns each contextual embedding $\mathbf{C}_i \in \mathbb{R}^{D'}$ to the most similar RVQ token $\mathbf{Q}_t^{(1)}$ embedding within a localized search window using cosine similarity. A dynamic window-shifting mechanism prevents alignment overlap and ensures sequence-wide consistency. If multiple RVQ tokens within the window share the highest similarity, \mathbf{C}_i is assigned to all corresponding positions in the aligned output \mathbf{C}^* , accounting for cases where a single linguistic token spans multiple speech frames. The resulting sequence $\mathbf{C}^* \in \mathbb{R}^{T' \times D'}$ enables timestep-level supervision. The full procedure is detailed in Algorithm 1.

> Algorithm 1 Window-Based Token Alignment **Require:** Contextual embeddings $\{C_i\}_{i=1}^n$, RVQ tokens $\{\mathbf{Q}_t^{(1)}\}_{t=1}^{T'}$, optional window size w 1: if w not provided then $w \leftarrow |T'/n|$ 2: 3: end if 4: Initialize aligned output $\mathbf{C}^* \in \mathbb{R}^{T' \times D'} \leftarrow 0$ 5: Initialize $\ell \leftarrow 0$ {last matched index} 6: **for** i = 1 to n **do** if dynamic window then 7: $s \leftarrow \ell + 1$ if i > 1, else 0 {start index} 8: $e \leftarrow \min(s + w, T')$ {end index} 9: 10: else $s \leftarrow (i-1) \cdot w, e \leftarrow \min(s+w, T')$ 11: end if 12: 13: Compute cosine similarity $\alpha_t = \cos(\mathbf{C}_i, \mathbf{Q}_t^{(1)}) \text{ for } t \in [s, e)$ Let $\tau \leftarrow \max_t \alpha_t$ {maximum similarity} 14: $\mathcal{T}_i \leftarrow \{t \mid \alpha_t \ge \tau\}$ 15: for each $t \in \mathcal{T}_i$ do 16: $\mathbf{C}_t^* \leftarrow \mathbf{C}_i$ 17: end for 18: $\ell \leftarrow \max(\mathcal{T}_i)$ 19: 20: end for 21: return C*

> Finally, we apply temporally aligned contextual supervision using a timestep-level distillation loss:

$$\mathcal{L}_{\text{distill}} = -\frac{1}{T'} \sum_{t=1}^{T'} \log \sigma \left(\cos \left(\mathbf{Q}_t^{\prime(1)}, \mathbf{C}_t^* \right) \right), \quad (4)$$

where $\mathbf{Q}^{\prime(1)} = \mathbf{Q}^{(1)}\mathbf{W} \in \mathbb{R}^{T' \times D'}$ is the linearly projected RVQ output, and $\sigma(\cdot)$ denotes the sigmoid function. This loss enforces temporally precise alignment between acoustic tokens and their corresponding contextual representations, encouraging modality-aware token learning.

358

359

360

361

362

363

364

366

367

368

369

370

371

372

373

374

376

377

378

379

381

382

383

384

386

389

390

391

392

393

394

395

396

397

398

400

401

402

403

404

2.4 Architecture and Training Objective

We build on widely adopted neural codec architectures and training objectives, following (Défossez et al., 2022; Zhang et al., 2024; Ahasan et al., 2024), to establish a strong and reliable foundation. We contribute to enhancing the learned representations through semantic and contextual supervision and fusion without altering the model architecture.

We use wav2vec 2.0 (base-960h) as the ASR model A (Baevski et al., 2020), BERT (bert-baseuncased) as the language model B (Devlin et al., 2019a), and HuBERT (base-ls960) as the selfsupervised speech model H (Hsu et al., 2021). All pre-trained models are frozen during training. The speech tokenizer consists of an encoder E, an RVQ module with 8 quantization layers (codebooks) of size 1024, a decoder D, and three discriminators (multi-period, multi-scale, and multi-scale STFT). Architectural details are provided in Appendix B.1.

Quantization operates on 50 Hz frame rates. The encoder and RVQ use an embedding dimension of D = 1024, while the pre-trained langauge and speech model have D' = 768. Cross-Attentions are implemented using 8-heads. The dropout masks \mathcal{D}_S and \mathcal{D}_C are applied at a rate of 10%.

Training Objective. We also adopt a multiobjective training setup grounded in established neural codec practices. This includes time-domain reconstruction loss \mathcal{L}_{time} , frequency-domain reconstruction loss \mathcal{L}_{freq} , adversarial loss \mathcal{L}_{gen} , feature matching loss \mathcal{L}_{feat} , and RVQ commitment loss \mathcal{L}_{commit} (see Appendix B.2 for details).

To further enhance representation quality, we introduce two auxiliary supervision objectives: a global distillation loss as $\mathcal{L}_{distill}$ (Sec. 2.3.2) and a temporally aligned contextual loss as $\mathcal{L}_{distill}$ (Sec. 2.3.3). $\mathcal{L}_{distill}$ is set to 0, when applying the Latent Representation Fusion (Sec. 2.3.1).

The final training objective is a weighted sum:

$$\mathcal{L}_{\text{total}} = \lambda_{\text{time}} \mathcal{L}_{\text{time}} + \lambda_{\text{freq}} \mathcal{L}_{\text{freq}} + \lambda_{\text{gen}} \mathcal{L}_{\text{gen}} + \lambda_{\text{feat}} \mathcal{L}_{\text{feat}} + \lambda_{\text{commit}} \mathcal{L}_{\text{commit}}$$
(5)
+ $(\lambda_{\text{distill}} \mathcal{L}_{\text{distill}} \text{ or } 0)$

.....

355

336

337

340

341

343

344

345

346

347

349

351

353

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440 441

442

443

444

445

446

447

448

449

450

451

452

2.5 Downstream Extension to TTS Model

We extend the learned discrete token representations to a downstream text-to-speech (TTS) task, following the neural codec language modeling framework and objective used in prior work (Wang et al., 2023; Zhang et al., 2024; Ahasan et al., 2024). In this paradigm, speech synthesis is performed by predicting quantized acoustic tokens produced by the RVQ and decoded by a neural codec.

For this, we propose FuseCodec-TTS, an extension of FuseCodec-Fusion trained with either Latent Representation Fusion (see §2.3.1) or FuseCodec-Distill using Global Semantic-Contextual Supervision (see §2.3.2). This allows the TTS model to operate on discrete speech tokens enriched with semantic and contextual information.

Given a phoneme sequence p and an acoustic prompt $\mathbf{A} \in \mathbb{R}^{\tau \times K}$ extracted from a reference utterance using FuseCodec, the goal is to predict a sequence of discrete token indices $q^{(1)}, \ldots, q^{(K)}$, corresponding to the K RVQ layers.

To model coarse content and prosodic structure, we autoregressively predict the token indices $q^{(1)}$ from the first quantizer using a decoder-only Transformer conditioned on the phoneme sequence **p**. The autoregressive (AR) training objective is:

$$\mathcal{L}_{AR} = -\log \prod_{i=1}^{T'} p(q_i^{(1)} \mid q_{$$

To capture fine-grained acoustic details, we use a non-autoregressive model to predict $q^{(k)}$ for each k = 2, ..., K, conditioned on the previously predicted layers $q^{(<k)}$, the phoneme sequence **p**, and the acoustic prompt **A**. The non-autoregressive (NAR) training objective is:

$$\mathcal{L}_{\text{NAR}} = -\log \prod_{k=2}^{K} p(q^{(k)} \mid q^{((7)$$

Both AR and NAR token generators are implemented using 12-layer Transformers with 16 attention heads, 1024-dimensional embeddings, 4096dimensional feed-forward layers, and a dropout rate of 0.1. The predicted token indices are mapped to their corresponding quantized embeddings $\mathbf{Q}^{(k)}$, which are then passed to FuseCodec's decoder to reconstruct the synthesized speech waveform.

3 Experimental Setup

Dataset. Following prior work in speech tokenization (Zhang et al., 2024; Ahasan et al., 2024), we train FuseCodec on the LibriSpeech (Panayotov et al., 2015) train-clean-100 subset, which contains 100 hours of English speech from 251 speakers, sampled at 16 kHz. During training, we randomly crop 3-second audio segments and reserve 100 samples for validation. For FuseCodec-TTS, we combine the train and dev subsets of LibriTTS (Zen et al., 2019), comprising 570 hours of speech. 453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

For evaluating FuseCodec, we use the LibriSpeech test-clean subset, which comprises 2,620 utterances held out entirely from training. This setup follows prior baselines (Zhang et al., 2024; Ahasan et al., 2024), though we evaluate on the full set rather than a sampled subset. For FuseCodec-TTS, we adopt two established benchmark protocols. In the LibriSpeech evaluation, following Wang et al. (2023), we select utterances between 4 and 10 seconds, yielding a 2.2-hour subset. For each synthesis, a 3-second enrollment segment is randomly cropped from a different utterance by the same speaker. In the VCTK evaluation, following Zhang et al. (2024), a 3-second prompt is selected or cropped from one utterance, and the transcript of a separate utterance from the same speaker serves as the synthesis target.

Training. FuseCodec is trained for 100 epochs on two A40 GPUs with a batch size of 6, using the Adam optimizer with a learning rate of 1×10^{-4} and exponential decay factor 0.98. FuseCodec-TTS is trained on A100 and L40S GPUs. The AR model is trained for 200 epochs, and the NAR model for 150 epochs. Training employs dynamic batching, with each batch containing up to 550 seconds of audio for AR and 100–200 seconds for NAR. We use the ScaledAdam optimizer with a learning rate of 5×10^{-2} and 200 warm-up steps.

Reproducibility. We provide a fully reproducible setup, including a Dockerized environment, source code, model checkpoints, and configuration files (anonymized for review; see Appendix A.)

Baselines. We compare FuseCodec against both established and recent strong baseline speech tokenizers, including EnCodec (24 kHz) (Défossez et al., 2022) and SpeechTokenizer (Zhang et al., 2024), as well as BigCodec (Xin et al., 2024), DAC (16 kHz) (Kumar et al., 2023), DM-Codec (LM+SM) (Ahasan et al., 2024) FACodec (Natural-Speech 3) (Ju et al., 2024), and Moshi (Défossez et al., 2024). All baseline results are obtained using official released checkpoints. For FuseCodec-TTS, we compare with neural codec language models that incorporate external representation guidance. Specifically, we compare against USLM

Table 1: Speech reconstruction results across content preservation and naturalness metrics. Orange and light orange cells indicate the best and second-best scores, respectively. Results show that FuseCodec variants outperform baselines by unifying contextual and semantic signals in the discrete speech representations.

Model	Content Preservation			Speech Naturalness					
110001	WER↓	WIL↓	STOI↑	ViSQOL↑	PESQ↑	UTMOS↑	Similarity↑		
BigCodec	4.58	7.45	0.93	3.02	2.68	3.44	0.996		
DAC	4.09	6.54	0.94	3.36	2.72	3.33	0.996		
DM-Codec	4.09	6.75	0.93	3.20	2.77	3.45	0.994		
EnCodec	4.04	6.58	0.92	3.06	2.31	2.41	0.980		
FACodec	4.11	6.58	0.95	3.11	2.89	3.45	0.996		
Mimi	11.61	18.05	0.85	2.49	1.69	2.28	0.934		
SpeechTokenizer	4.16	6.71	0.92	3.08	2.60	3.41	0.996		
FuseCodec (Baseline)	4.62	7.44	0.93	2.95	2.54	3.18	0.990		
FuseCodec-ContextAlign	4.15	6.70	0.93	3.18	2.85	3.65	0.995		
FuseCodec-Distill	4.09	6.60	0.94	3.43	3.06	3.65	0.996		
FuseCodec-Fusion	3.99	6.45	0.95	3.47	3.13	3.63	0.995		

505(from SpeechTokenizer) (Zhang et al., 2024) and506DM-Codec-TTS (Ahasan et al., 2024), using their507official released LibriTTS trained checkpoints.

508 Metrics. We evaluate FuseCodec using two complementary categories of metrics: Content Preservation and Speech Naturalness. To assess Content 510 Preservation, we transcribe generated speech using Whisper (medium) (Radford et al., 2023) and com-512 pare it to ground-truth text. We report Word Error 513 *Rate (WER)*, defined as WER = $\frac{S+D+I}{N}$, where S, 514 D, and I denote the number of substitutions, dele-515 tions, and insertions, and N is the number of words 516 in the reference. We also report Word Information 517 Lost (WIL), given by WIL = $1 - \frac{C}{N} + \frac{C}{P}$, where 518 C is the number of correct words, N is the num-519 ber of words in the reference, and P is the number of words in the prediction. Additionally, we in-521 clude Short-Time Objective Intelligibility (STOI), a 522 reference-based metric estimating intelligibility via short-time spectral similarity. For Speech Natural-524 ness, we evaluate perceptual and acoustic fidelity 525 using both reference-based and learned metrics. ViSQOL and PESQ assess perceptual quality by 527 modeling auditory similarity and signal distortion, respectively. We also report UTMOS for estimat-529 ing human-judged naturalness, which is a neural MOS predictor trained on large-scale human rat-531 ings. Lastly, we compute Similarity as the cosine similarity between L2-normalized speaker embed-533 dings extracted using WavLM-TDNN (Chen et al., 534 2022), reflecting speaker or content consistency. For FuseCodec-TTS, we omit metrics requiring 536 reference audio (e.g., STOI, ViSQOL, PESQ), as 537 exact references are unavailable in synthesis.

4 Experimental Results and Discussion

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

In this section, we evaluate our proposed methods on speech reconstruction quality (§4.1), their extension to speech synthesis (§4.2), and validate the contribution of each component through comprehensive ablation studies (Appendix C)).

4.1 Speech Reconstruction Evaluation

We evaluate our three proposed methods: (i) Latent Representation Fusion: FuseCodec-Fusion (Sec. 2.3.1) (ii) Global Semantic-Contextual Supervision: FuseCodec-Distill (Sec. 2.3.2), and (iii) Temporally Aligned Contextual Supervision: FuseCodec-ContextAlign (Sec. 2.3.3).

Results. The results in Table 1 show that FuseCodec improves performance across all metrics related to content preservation and speech naturalness. FuseCodec-Fusion performs best overall, achieving the lowest WER (3.99) and WIL (6.45), along with the highest STOI (0.95), reducing transcription error and improving intelligibility. It also achieves the highest scores in ViSQOL (3.47) and PESQ (3.13), reflecting superior perceptual quality. FuseCodec-Distill attains top scores in UT-MOS (3.65) and Similarity (0.996), while also ranking second in STOI (0.94), ViSQOL (3.43), and PESQ (3.06), demonstrating strong naturalness and speaker consistency. FuseCodec-ContextAlign also performs competitively, particularly in UTMOS (3.65) and Similarity (0.995), while showing consistent improvements over FuseCodec (Baseline).

Discussion. FuseCodec-Fusion achieves the best overall performance. Compared to EnCodec, which focuses purely on acoustic representations, and FACodec, which separates attribute learning without unifying representations, FuseCodec-

Table 2: Zero-shot TTS evaluation on LibriSpeech and VCTK. FuseCodec-TTS variants are compared to official neural codec-based TTS checkpoints trained on LibriTTS. Bold and underline indicate best and second-best scores. FuseCodec-TTS improves intelligibility, similarity, and naturalness via semantic-contextual aware tokenization.

	WER	Ļ	WIL	Ļ	Similari	ty↑	UTMOS ↑	
Model	LibriSpeech	VCTK	LibriSpeech	VCTK	LibriSpeech	VCTK	LibriSpeech	VCTK
DM-Codec-TTS	10.26	5.02	13.79	8.21	0.82	0.79	3.70	3.86
USLM	16.72	14.79	25.65	23.24	0.80	0.78	2.93	3.01
FuseCodec-Distill-TTS	8.55	3.66	12.07	6.02	0.82	0.78	3.55	3.75
FuseCodec-Fusion-TTS	<u>9.67</u>	<u>4.07</u>	<u>13.23</u>	7.18	0.83	0.79	<u>3.63</u>	<u>3.82</u>

Fusion incorporates both semantic and contextual 574 signals directly into the encoder's latent space. This enables the quantizer to learn a unified represen-576 tation aligned with both linguistic meaning and 577 acoustic structure. It also outperforms models like 578 DAC and BigCodec, which prioritize compression but lack representational alignment. FuseCodec-Distill improves upon SpeechTokenizer and Mimi, which distill only semantic representations from speech models and underperform on intelligibility 584 and quality. In contrast, FuseCodec-Distill supervises the quantized space with global contextual and semantic signals, promoting alignment with 586 high-level linguistic and acoustic content.

590

591

597

608

610

611

612

614

FuseCodec-ContextAlign introduces finegrained supervision by aligning discrete tokens with temporally matched contextual tokens, encouraging each token to reflect local linguistic context. Although its constrained alignment limits global contextual guidance, leading to slightly lower performance than FuseCodec-Fusion and FuseCodec-Distill, it still outperforms DM-Codec, improving intelligibility and speaker similarity. Overall, contextual guidance is most effective for content preservation, while semantic supervision enhances speech naturalness. FuseCodec-Fusion delivers the best balance, FuseCodec-Distill excels in speaker fidelity, and FuseCodec-ContextAlign offers interpretable gains. These results underscore the benefit of unifying multimodal representations.

4.2 Speech Synthesis Evaluation

We extend our methods to the zero-shot text-tospeech task and evaluate against neural codec TTS models that incorporate representational supervision for comparison. We deicde to adapt FuseCodec-Fusion to FuseCodec-Fusion-TTS and FuseCodec-Distill to FuseCodec-Distill-TTS, based on their superior performance.

Results. Table 2 shows that both FuseCodec-TTS variants outperform prior methods across most metrics on LibriSpeech and VCTK. FuseCodecDistill-TTS achieves the best content preservation, with the lowest WER (8.55/3.66) and WIL (12.07/ 6.02), surpassing both DM-Codec-TTS and USLM. FuseCodec-Fusion-TTS delivers the highest perceptual quality, achieving the top speaker similarity (0.83 / 0.79), while also maintaining strong intelligibility with the second-best UTMOS (3.63 / 3.82), WER (9.67 / 4.07), and WIL (13.23 / 7.18). 615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

Discussion. FuseCodec-Fusion-TTS leads in perceptual quality and speaker similarity. Unlike DM-Codec-TTS, which lacks precise alignment, and USLM, which incorporates only semantic features, FuseCodec-Fusion-TTS integrates both semantic and contextual signals directly into the encoder's latent space. This allows the quantizer to capture expressive prosody and speaker identity, resulting in more natural and coherent speech. FuseCodec-Distill-TTS achieves the highest intelligibility and transcription accuracy. In contrast to USLM's lack of contextual grounding and DM-Codec-TTS's limited supervision, it distills global semantic-contextual representations into the quantized token space, enhancing alignment with semantic and contextual info. While FuseCodec-Fusion-TTS excels in naturalness and speaker fidelity, FuseCodec-Distill-TTS offers stronger linguistic precision. This trade-off reflects the complementary strengths of each variant and underscores the importance of integrating semantic-contextual fusion or supervision into speech tokenization.

5 Conclusion

We introduced FuseCodec, a unified speech tokenization framework that integrates acoustic, semantic, and contextual signals via multimodal representation fusion and supervision. Our methods enable fine-grained alignment and achieve state-ofthe-art results on speech reconstruction, improving intelligibility, quality, and speaker similarity. These findings highlight the value of semantic and contextual grounding in discrete speech modeling.

Limitations

References

While FuseCodec demonstrates strong perfor-

mance in incorporating semantic and contextual

representation into discrete speech tokenization,

our approach relies on frozen pretrained language

and speech models, which may limit adaptability

to non-English or low-resource languages. More-

over, although we focus on zero inference over-

head, training requires additional computational re-

sources due to multimodal guidance. Future work

could explore lightweight, adaptive guidance mech-

Md Mubtasim Ahasan, Md Fahim, Tasnim Mohiud-

din, A K M Mahbubur Rahman, Aman Chadha,

Tariq Iqbal, M Ashraful Amin, Md Mofijul Islam,

and Amin Ahsan Ali. 2024. Dm-codec: Distilling multimodal representations for speech tokenization.

Alexei Baevski, Henry Zhou, Abdelrahman Mohamed,

Yoshua Bengio, Nicholas Léonard, and Aaron Courville.

Zalán Borsos, Raphaël Marinier, Damien Vincent, Eu-

gene Kharitonov, Olivier Pietquin, Matt Sharifi,

Dominik Roblek, Olivier Teboul, David Grangier,

Marco Tagliasacchi, and Neil Zeghidour. 2023. Au-

diolm: A language modeling approach to audio gen-

eration. IEEE/ACM Trans. Audio, Speech and Lang.

Annemarie C. Brown, Eva Childers, Elijah F. W. Bowen,

Gabriel A. Zuckerberg, and Richard Granger. 2022.

Phonemes in continuous speech are better recognized

in context than in isolation. Frontiers in Communica-

Sanyuan Chen, Chengyi Wang, Zhengyang Chen,

Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki

Kanda, Takuya Yoshioka, Xiong Xiao, Jian Wu, Long

Zhou, Shuo Ren, Yanmin Qian, Yao Qian, Jian Wu,

Michael Zeng, Xiangzhan Yu, and Furu Wei. 2022.

Wavlm: Large-scale self-supervised pre-training for

full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing*, 16(6):1505–1518.

Djork-Arné Clevert, Thomas Unterthiner, and Sepp

Hochreiter. 2016. Fast and accurate deep network

learning by exponential linear units (elus). Preprint,

2013. Estimating or propagating gradients through stochastic neurons for conditional computation.

and Michael Auli. 2020. wav2vec 2.0: A framework

for self-supervised learning of speech representations.

anisms and broader language coverage.

Preprint, arXiv:2410.15017.

Preprint, arXiv:2006.11477.

Preprint, arXiv:1308.3432.

Proc., 31:2523-2533.

tion, Volume 7 - 2022.

arXiv:1511.07289.

657

66 66

66

660

-
- 66
- 67 67
- 67

674

- 67
- 67
- 681
- 68
- 68 68

68

68

69

69 69

- 69
- 695 696

6

700 701

- 7
- 703 704 705

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019a. Bert: Pre-training of deep bidirectional transformers for language understanding. *Preprint*, arXiv:1810.04805. 706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019b. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi. 2022. High fidelity neural audio compression. *Preprint*, arXiv:2210.13438.
- Alexandre Défossez, Laurent Mazaré, Manu Orsini, Amélie Royer, Patrick Pérez, Hervé Jégou, Edouard Grave, and Neil Zeghidour. 2024. Moshi: a speech-text foundation model for real-time dialogue. *Preprint*, arXiv:2410.00037.
- Mark Hallap, Emmanuel Dupoux, and Ewan Dunbar. 2023. Evaluating context-invariance in unsupervised speech representations. *Preprint*, arXiv:2210.15775.
- Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM transactions on audio*, *speech, and language processing*, 29:3451–3460.
- Ahmed Hussen Abdelaziz, Barry-John Theobald, Paul Dixon, Reinhard Knothe, Nicholas Apostoloff, and Sachin Kajareker. 2020. Modality dropout for improved performance-driven talking faces. In Proceedings of the 2020 International Conference on Multimodal Interaction, ICMI '20, page 378–386, New York, NY, USA. Association for Computing Machinery.
- Zeqian Ju, Yuancheng Wang, Kai Shen, Xu Tan, Detai Xin, Dongchao Yang, Yanqing Liu, Yichong Leng, Kaitao Song, Siliang Tang, Zhizheng Wu, Tao Qin, Xiang-Yang Li, Wei Ye, Shikun Zhang, Jiang Bian, Lei He, Jinyu Li, and Sheng Zhao. 2024. Naturalspeech 3: Zero-shot speech synthesis with factorized codec and diffusion models. *Preprint*, arXiv:2403.03100.
- Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. 2020. Hifi-gan: generative adversarial networks for efficient and high fidelity speech synthesis. In *Proceedings* of the 34th International Conference on Neural Information Processing Systems, NIPS '20, Red Hook, NY, USA. Curran Associates Inc.
- Kundan Kumar, Rithesh Kumar, Thibault de Boissiere, Lucas Gestin, Wei Zhen Teoh, Jose Sotelo, Alexandre de Brebisson, Yoshua Bengio, and Aaron Courville. 2019. Melgan: Generative adversarial networks for conditional waveform synthesis. *Preprint*, arXiv:1910.06711.

- 765 766 767
- 76
- 77
- 772 773 774 775
- 776 777 778 779 780
- 7
- 784 785 786 786
- 788 789 790 791
- 792 793 794 795
- 796
- 797 798
- 79 80 80
- 80 80 80
- 806 807
- 808 809
- 810 811
- 812 813

815

- 816 817
- 817 818

819

- Rithesh Kumar, Prem Seetharaman, Alejandro Luebs, Ishaan Kumar, and Kundan Kumar. 2023. Highfidelity audio compression with improved rvqgan. *Preprint*, arXiv:2306.06546.
- Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. Librispeech: An asr corpus based on public domain audio books. In 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 5206–5210.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.
 - Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine Mcleavey, and Ilya Sutskever. 2023.
 Robust speech recognition via large-scale weak supervision. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 28492–28518. PMLR.
- Craig W Schmidt, Varshini Reddy, Haoran Zhang, Alec Alameddine, Omri Uzan, Yuval Pinter, and Chris Tanner. 2024. Tokenization is more than compression. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 678–702, Miami, Florida, USA. Association for Computational Linguistics.
- Arnon Turetzky and Yossi Adi. 2024. Last: Language model aware speech tokenization. *Preprint*, arXiv:2409.03701.
- Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. 2018. Neural discrete representation learning. *Preprint*, arXiv:1711.00937.
- Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang, Long Zhou, Shujie Liu, Zhuo Chen, Yanqing Liu, Huaming Wang, Jinyu Li, Lei He, Sheng Zhao, and Furu Wei. 2023. Neural codec language models are zero-shot text to speech synthesizers. *Preprint*, arXiv:2301.02111.
- Detai Xin, Xu Tan, Shinnosuke Takamichi, and Hiroshi Saruwatari. 2024. Bigcodec: Pushing the limits of low-bitrate neural speech codec. *Preprint*, arXiv:2409.05377.
- Dongchao Yang, Songxiang Liu, Rongjie Huang, Jinchuan Tian, Chao Weng, and Yuexian Zou.
 2023. Hifi-codec: Group-residual vector quantization for high fidelity audio codec. *Preprint*, arXiv:2305.02765.
- Zhen Ye, Peiwen Sun, Jiahe Lei, Hongzhan Lin, Xu Tan, Zheqi Dai, Qiuqiang Kong, Jianyi Chen, Jiahao Pan, Qifeng Liu, Yike Guo, and Wei Xue. 2024.

Codec does matter: Exploring the semantic shortcoming of codec for audio language model. *Preprint*, arXiv:2408.17175. 820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

- Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi. 2022. Soundstream: An end-to-end neural audio codec. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 30:495–507.
- Heiga Zen, Viet Dang, Rob Clark, Yu Zhang, Ron J. Weiss, Ye Jia, Zhifeng Chen, and Yonghui Wu. 2019. Libritts: A corpus derived from librispeech for textto-speech. *Preprint*, arXiv:1904.02882.
- Xin Zhang, Dong Zhang, Shimin Li, Yaqian Zhou, and Xipeng Qiu. 2024. Speechtokenizer: Unified speech tokenizer for speech language models. In *The Twelfth International Conference on Learning Representations*.

- 000 839
- 840
- 841
- 842
- 0-
- 04 84
- 8/
- 8
- 8

851

855

857

864

871

873

874

876

881

A Resources

We provide all necessary resources to ensure full reproducibility of our models and results. All links are anonymized for double-blind peer review.

- **Docker:** A containerized environment with all required Python packages for training. LINK
- Code and Configuration: Full codebase for preprocessing, training, and inference. LINK
- Model Checkpoints: Trained model weights. LINK

B Tokenizer Design and Loss Functions

B.1 Model Details

To implement a strong speech tokenizer baseline, we adopt a standard neural codec architecture and discriminator setup commonly used in prior work (Défossez et al., 2022; Zeghidour et al., 2022).

Encoder and Decoder. The Encoder consists of an initial 1D convolutional layer with 32 channels and a kernel size of 7, followed by 4 stacked residual blocks. Each block includes two dilated convolutions with a (3, 1) kernel and no dilation expansion (dilation = 1), a residual connection, and a strided convolutional layer for temporal downsampling. Stride values across the blocks are set to 2, 4, 5, and 8, with kernel sizes for the downsampling layers set to twice the corresponding stride. Channel dimensions double at each downsampling stage. The encoder then includes a two-layer BiLSTM, and concludes with a 1D convolution (kernel size 7) to project to the target embedding dimension. ELU (Clevert et al., 2016) is used as the activation function, and layer normalization or weight normalization is applied depending on the layer. The Decoder mirrors the encoder architecture, with the only difference being the use of transposed convolutions in place of strided convolutions to reverse the downsampling steps, and the inclusion of LSTM layers to restore temporal resolution.

Residual Vector Quantizer. The Residual Vector Quantizer (RVQ) module discretizes the encoder's continuous latent representations into a sequence of codebook indices. Specifically, we quantize the encoder latent tensor of shape [B, D, T] using 8 residual codebooks, each with 1024 codebook entries. Each subsequent codebook quantizes the residual error of the previous one. Codebook entries are updated using an exponential moving

average with a decay factor of 0.99. To prevent codebook collapse, unused entries are randomly resampled using vectors from the current batch. The RVQ output is a discrete tensor of shape $[B, N_q, T]$, where N_q is the number of active quantizers. The indices are mapped back to the original latent space by summing the corresponding codebook embeddings and are then fed into the decoder to reconstruct the input. A straight-through estimator (Bengio et al., 2013) is used to propagate gradients through the quantizer. 885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

Discriminators. We utilize discriminators to guide the generators (Encoder, RVQ, and Decoder) to reconstruct speech more closely to the original. We make use of three distinct discriminators: a Multi-Scale STFT (MS-STFT) discriminator, a Multi-Scale Discriminator (MSD), and a Multi-Period Discriminator (MPD). The MS-STFT discriminator, proposed by (Défossez et al., 2022), works on multiple resolutions of the complexvalued short-time Fourier transform (STFT). It treats the real and imaginary parts as concatenated and applies a sequence of 2D convolutional layers. The initial layer uses a kernel size of 3×8 with 32 channels. This is followed by convolutions with increasing temporal dilation rates (1, 2,and 4) and a stride of 2 along the frequency axis. A final 3×3 convolution with stride 1 outputs the discriminator prediction. The MSD processes the raw waveform at various temporal scales using progressively downsampled versions of the input. We adopt the configuration from (Zeghidour et al., 2022), which was originally based on (Kumar et al., 2019). Similarly, the MPD, introduced by (Kong et al., 2020), models periodic structure in the waveform by reshaping it into a 2D input with unique periodic patterns. For consistency, we standardize the number of channels in both the MSD and MPD to match those in the MS-STFT discriminator.

B.2 Training Objective

To ensure that FuseCodec learns discrete speech representations, we ground our training objective on proven techniques, following (Défossez et al., 2022; Zhang et al., 2024; Ahasan et al., 2024).

Reconstruction loss. Let x and \hat{x} denote the original and reconstructed speech waveforms, respectively. For spectral comparisons, we define 64-bin Mel-spectrograms $\mathbf{M}_i(\cdot)$ using STFTs with window size 2^i and hop size $2^i/4$, where $i \in \mathcal{E} = \{5, \ldots, 11\}$ indexes different resolution scales. We compute the time-domain $\mathcal{L}_{\text{time}}$ and frequency-

$$\mathcal{L}_{ ext{freq}} = \sum_{i \in \mathcal{S}} igg(\| \mathbf{M}_i(\mathbf{x}) - \mathbf{M}_i(\hat{\mathbf{x}}) \|_1$$

 $\mathcal{L}_{time} = \|\mathbf{x} - \hat{\mathbf{x}}\|_1$

944

945

948

951

954

955

957

961

962

963

$$+\|\mathbf{M}_{i}(\mathbf{x}) - \mathbf{M}_{i}(\hat{\mathbf{x}})\|_{2}\right)$$
(9)

Adversarial loss. To reduce the discriminability of reconstructed speech, we adopt a GANbased training objective with a set of discriminators $\{D^{(i)}\}_{i=1}^{d}$, including multi-period (MPD), multiscale (MSD), and multi-scale STFT (MS-STFT) variants (see Appendix B for details). The generator \mathcal{L}_{gen} and discriminator \mathcal{L}_{disc} losses are computed as:

$$\mathcal{L}_{\text{gen}} = \frac{1}{d} \sum_{i=1}^{d} \max\left(0, \ 1 - D^{(i)}(\hat{\mathbf{x}})\right) \qquad (10)$$

$$\mathcal{L}_{\text{disc}} = \frac{1}{d} \sum_{i=1}^{d} \left[\max(0, \ 1 - D^{(i)}(\mathbf{x})) + \max(0, \ 1 + D^{(i)}(\hat{\mathbf{x}})) \right]$$
(11)

Let $D_j^{(i)}(\cdot)$ denote the output of the *j*-th layer of $D^{(i)}$, with ℓ total layers. We include a feature $\mathcal{L}_{\text{feat}}$ matching loss to stabilize training and align intermediate features as:

$$\mathcal{L}_{\text{feat}} = \frac{1}{d\ell} \sum_{i=1}^{d} \sum_{j=1}^{\ell} \frac{\|D_j^{(i)}(\mathbf{x}) - D_j^{(i)}(\hat{\mathbf{x}})\|_1}{\text{mean}\left(\|D_j^{(i)}(\mathbf{x})\|_1\right)} \quad (12)$$

Commitment Loss. To ensure encoder outputs align closely with their quantized representations, we apply a commitment penalty during residual vector quantization (RVQ). Let \mathbf{r}_j denote the residual vector at step $j \in \{1, ..., q\}$, and \mathbf{c}_j be its corresponding nearest codebook entry, we calculate commitment loss $\mathcal{L}_{\text{commit}}$ as:

964
$$\mathcal{L}_{\text{commit}} = \sum_{j=1}^{q} \|\mathbf{r}_j - \mathbf{c}_j\|_2^2 \qquad (13)$$

C Ablation Studies

(8)

We ablate and investigate each design choice and the necessity of components in our proposed methodology for FuseCodec. All model hyperparameters, training procedures, and configurations are kept fixed, except for the specific changes introduced in each ablation setup. 965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

C.1 Ablation: Attention-Projection Configuration in Representation Fusion

Setup. We investigate the impact of changing the attention-projection configuration in FuseCodec-Fusion (Section 2.3.1). The selected method, **Cross-Before**, applies multi-head cross-attention prior to projection:

$$\mathbf{S}' = \operatorname{CrossAttention}(\tilde{\mathbf{S}}, \tilde{\mathbf{C}}, \tilde{\mathbf{C}}) \mathbf{W}_{S},$$

$$\mathbf{C}' = \operatorname{CrossAttention}(\tilde{\mathbf{C}}, \tilde{\mathbf{S}}, \tilde{\mathbf{S}}) \mathbf{W}_{C},$$
 (14)

where $\tilde{\mathbf{S}}, \tilde{\mathbf{C}} \in \mathbb{R}^{T' \times D'}$ are broadcasted global semantic and contextual vectors. We compare this with the following ablated variants:

None, which skips attention and directly applies projection:

$$\mathbf{S}' = \mathbf{S}\mathbf{W}_S,$$

$$\mathbf{C}' = \tilde{\mathbf{C}}\mathbf{W}_C$$
(15)

Self-Before, which applies self-attention before projection:

$$\mathbf{S}' = \text{SelfAttention}(\tilde{\mathbf{S}}, \tilde{\mathbf{S}}, \tilde{\mathbf{S}}) \mathbf{W}_S,$$

$$\mathbf{C}' = \text{SelfAttention}(\tilde{\mathbf{C}}, \tilde{\mathbf{C}}, \tilde{\mathbf{C}}) \mathbf{W}_C$$
 (16)

Self-After, which projects first and then applies self-attention:

$$\mathbf{S}' = \text{SelfAttention}(\tilde{\mathbf{S}}\mathbf{W}_S),$$

$$\mathbf{C}' = \text{SelfAttention}(\tilde{\mathbf{C}}\mathbf{W}_C)$$
(17)

Cross-After, which applies projection before crossattention:

$$\mathbf{S}' = \operatorname{CrossAttention}(\mathbf{\tilde{S}W}_S, \mathbf{\tilde{C}W}_C, \mathbf{\tilde{C}W}_C),$$
$$\mathbf{C}' = \operatorname{CrossAttention}(\mathbf{\tilde{C}W}_C, \mathbf{\tilde{S}W}_S, \mathbf{\tilde{S}W}_S)$$
(18)

Results. Table 3 shows the results of five vari-
ants. The selected Cross-Before setup achieves the
highest performance on intelligibility STOI (0.95),
and all naturalness metrics: ViSQOL (3.47), PESQ
(3.13), and second-best UTMOS (3.63). Self-
Before yields the best WER (3.92) and WIL (6.36),
and second-best ViSQOL (3.43), PESQ (3.05), and
1001

Table 3: Ablation of attention-projection configurations in multimodal latent fusion. **Cross** variants incorporate cross-modal attention between semantic and contextual signals, while **Self** variants apply self-attention. **Before** applies attention prior to projection into the encoder's latent space, whereas **After** applies attention post-projection. **None** uses direct projection without attention. *Applying cross-modal attention before projection consistently improves content preservation and speech naturalness by enabling richer multimodal interactions in the original dimension.*

Model Variant	Attn-Proi Type	Conte	nt Preser	vation	Speech Naturalness				
		WER↓	WIL↓	STOI↑	ViSQOL↑	PESQ↑	UTMOS↑	Similarity↑	
FuseCodec-Fusion	None	4.10	6.60	0.93	3.26	2.92	3.65	0.995	
FuseCodec-Fusion	Self-After	4.07	6.61	0.93	3.26	2.95	3.63	0.995	
FuseCodec-Fusion	Self-Before	3.92	6.36	0.94	3.43	<u>3.05</u>	3.59	0.995	
FuseCodec-Fusion	Cross-After	4.17	6.70	0.93	3.28	2.90	3.61	0.995	
FuseCodec-Fusion	Cross-Before	<u>3.99</u>	6.45	0.95	3.47	3.13	<u>3.63</u>	0.995	

STOI (0.94). The **None** and **Cross-After** configurations perform comparatively worse across intelligibility and naturalness.

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1014

1015

1018

1019

1020

1021

1022

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1038

Discussion. These results demonstrate that the configuration of attention relative to projection significantly impacts the effectiveness of representation fusion. The best-performing method, **Cross-Before**, applies cross-modal attention in the original lower-dimensional space. This enables richer semantic-contextual interactions to be captured before transformation into the higher-dimensional encoder space, leading to improved intelligibility and perceptual quality.

Self-Before performs competitively by achieving the best WER and WIL, suggesting that intramodal structuring of global feature representations also benefits the fusion approach. However, the absence of explicit cross-modal exchange limits its effectiveness on naturalness metrics such as UT-MOS and PESQ.

By contrast, **Cross-After** performs poorly, indicating that applying cross-attention after projection diminishes its effectiveness. Suggesting that once projected into the higher-dimensional space, the global vectors lose semantic coherence, resulting in less expressive fusion and lower audio quality.

Finally, removing attention (**None**) results in the weakest performance on intelligibility and perceptual scores, despite yielding the highest UTMOS. This indicates that even unstructured modality signals can enhance naturalness, but without alignment through attention mechanisms, they fail to deliver consistent semantic-contextual grounding.

Overall, these results confirm that performing attention prior to projection, especially cross-modal attention, is essential for extracting the most benefit from semantic-contextual signals during fusion.

C.2 Ablation: Attention-Guidance Configuration in Semantic-Contextual Guidance

1039

1040

1041

1042

1043

1044

1045

1046

1047

1049

1050

1051

1052

1053

1054

1058

1059

1061

Setup. We study the impact of attention configuration and guidance modality used in the distillation objective. Our method, FuseCodec-Distill, introduces timestep-aligned supervision using global contextual and semantic signals (Section 2.3.2). The selected configuration, **None + Semantic-Contextual**, projects the first-layer RVQ tokens $Q^{(1)}$ and computes cosine similarity with both semantic and contextual guidance vectors:

$$\mathcal{L}_{\text{distill}} = -\frac{1}{T'} \sum_{t=1}^{T'} \log \sigma \left(\frac{1}{2} \left[\cos \left(\mathbf{Q}_t^{\prime(1)}, \tilde{\mathbf{S}}_t \right) + \cos \left(\mathbf{Q}_t^{\prime(1)}, \tilde{\mathbf{C}}_t \right) \right] \right)$$
(19)

We compare this against three ablated variants: **None + Contextual**, which excludes both attention and semantic guidance:

$$\mathcal{L}_{\text{distill}} = -\frac{1}{T'} \sum_{t=1}^{T'} \log \sigma \left(\cos \left(\mathbf{Q}_t^{\prime(1)}, \tilde{\mathbf{C}}_t \right) \right)$$
(20)

Cross + Contextual, which introduces crossattention between contextual vectors and projected RVQ tokens:

$$\tilde{\mathbf{C}} = \text{CrossAttention}(\tilde{\mathbf{C}}, \mathbf{Q}^{\prime(1)}, \mathbf{Q}^{\prime(1)})$$
 (21)

Cross + Semantic-Contextual, which includes cross-attention but retains both guidance signals.

Results. Table 4 reports the performance across1062four configurations. The best-performing variant1063is None + Semantic-Contextual, achieving the1064lowest WER (4.09) and WIL (6.60), and highest1065

Table 4: Ablation of attention and guidance strategies in semantic-contextual distillation. **Cross** variants apply cross-attention between contextual embeddings and discrete tokens, while **None** applies supervision directly. **Semantic-Contextual** combines both global semantic and contextual signals. *Direct supervision using both signals achieves the best intelligibility and perceptual quality by preserving global structure*.

Model Variant	Attention	Guidance	Conte	nt Preser	vation	Speech Naturalness				
Woder variant / Attention		Guidanee	WER↓	WIL↓	STOI↑	ViSQOL↑	PESQ↑	UTMOS↑	Similarity↑	
FuseCodec-Distill	None	Contextual	4.20	6.77	0.93	3.13	2.74	3.60	0.995	
FuseCodec-Distill	Cross	Contextual	4.18	<u>6.75</u>	0.93	<u>3.21</u>	2.83	3.60	<u>0.995</u>	
FuseCodec-Distill	None	Semantic-Contextual	4.09	6.60	0.94	3.43	3.06	3.65	0.996	
FuseCodec-Distill	Cross	Semantic-Contextual	4.21	6.82	0.93	3.18	<u>2.84</u>	<u>3.62</u>	0.994	

scores on STOI (0.940), ViSQOL (3.43), PESQ (3.06), UTMOS (3.65), and Similarity (0.996). The second-best results are obtained by **Cross + Contextual**, but excluding semantic guidance or using attention degrades performance across all metrics.

1066

1067

1068

1072

1073

1074

1075

1078

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1094

1095

1096

1098

1099

1100

1101

1102

1103

Discussion. These results show that including both semantic and contextual supervision is essential for improving the quantization quality of the discrete tokens. The **None + Semantic-Contextual** configuration outperforms all others, highlighting that cosine-based alignment with both modalities provides the most stable and effective guidance during quantized representation learning.

Introducing cross-attention (**Cross**) reduces performance, suggesting that attention distorts the global nature of the guidance signals and makes supervision less consistent across time. The **Cross + Semantic-Contextual** variant also underperforms, despite having access to both guidance sources, indicating that attention interferes with their inherent structure and alignment function.

The **Contextual-only** variants perform comparatively worse, confirming that semantic signals play an important role in guiding the learned representations toward higher-level content fidelity and improved intelligibility.

Overall, these findings support using both guidance signals in their original global forms and applying them directly, without attention, to ensure stable, timestep-aligned distillation.

C.3 Ablation: Fixed vs. Dynamic Window Configuration in Temporal Alignment

Setup. We investigate the effect of fixed versus dynamic windowing in the token alignment algorithm (Algorithm 1). Our full method, FuseCodec-ContextAlign, aligns each contextual embedding $\mathbf{C}_i \in \mathbb{R}^{D'}$ to a localized region of RVQ tokens $\{\mathbf{Q}_t^{(1)}\}_{t=1}^{T'}$ based on cosine similarity. The selected

configuration, **Dynamic-window Contextual** (see Section 2.3.3), dynamically adjusts the alignment window for each C_i , using the index of the previous match to guide the next search range. This contentaware strategy produces a temporally aligned sequence $C^* \in \mathbb{R}^{T' \times D'}$, which is used to compute a timestep-level distillation loss:

$$\mathcal{L}_{\text{align}} = -\frac{1}{T'} \sum_{t=1}^{T'} \log \sigma \left(\cos \left(\mathbf{Q}_t^{\prime(1)}, \mathbf{C}_t^* \right) \right)$$
(22)

We compare this setup against the following ablated variants:

Fixed-window Contextual, which uses a fixed alignment window of size $w = \lfloor T'/n \rfloor$, where T'is the RVQ sequence length and n is the number of contextual embeddings. Each C_i is aligned to the most similar token $Q_t^{(1)}$ within its predefined window.

Fixed-window Semantic-Contextual, which adds semantic supervision using semantic representations $\{S_i\}_{i=1}^m$, in addition to contextual representations aligned via a fixed-window token alignment. Since both semantic and RVQ tokens are extracted at the same frame rate, they are inherently timealigned, requiring no additional alignment. The combined loss is:

$$\begin{aligned} \mathcal{L}_{\text{align}} &= -\frac{1}{T'} \sum_{t=1}^{T'} \log \sigma \left(\frac{1}{2} \left[\cos \left(\mathbf{Q}_t'^{(1)}, \mathbf{C}_t^* \right) \right. \\ &+ \cos \left(\mathbf{Q}_t'^{(1)}, \mathbf{S}_t \right) \right] \end{aligned}$$

(23)

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

Dynamic-window Semantic-Contextual, which replaces the fixed window with a dynamic alignment strategy, while also incorporating direct supervision from semantic embeddings $\{S_t\}$.

Results. As shown in Table 5, the **Dynamic-**1133window Contextual configuration achieves the1134

Table 5: Ablation of windowing and guidance strategies in temporally aligned contextual supervision. **Dynamic** variants adapt the alignment window per token based on content similarity, while **Fixed** variants use a uniform window. **Semantic-Contextual** combines semantic and contextual signals for supervision. *Dynamic windowing consistently improves intelligibility and clarity by enabling finer temporal alignment of contextual embeddings.*

Model Variant	Window	Guidance	Content Preservation			Speech Naturalness			
			WER↓	WIL↓	STOI↑	ViSQOL↑	PESQ↑	UTMOS↑	Similarity↑
FuseCodec-ContextAlign FuseCodec-ContextAlign FuseCodec-ContextAlign FuseCodec-ContextAlign	Fixed Dynamic Fixed Dynamic	Contextual Contextual Semantic-Contextual Semantic-Contextual	4.26 4.15 4.30 <u>4.21</u>	6.88 6.70 6.88 <u>6.78</u>	0.92 0.93 0.92 <u>0.93</u>	3.19 <u>3.18</u> 3.10 3.12	2.71 2.85 2.62 <u>2.72</u>	3.58 3.65 <u>3.74</u> 3.75	0.994 0.995 0.995 0.995

best performance across content preservation met-1135 rics, achieving the lowest WER (4.15), WIL (6.70), 1136 and highest STOI (0.93). It also performs strongly 1137 in terms of speech naturalness, with the best PESQ 1138 (2.85), high ViSQOL (3.18), and top Similarity 1139 (0.995). The Dynamic Semantic-Contextual vari-1140 ant achieves the best UTMOS (3.75), second-best 1141 WER (4.21) and WIL (6.78), and matches the top 1142 Similarity. By contrast, both **Fixed-window** config-1143 urations obtains lower scores across most metrics, 1144 particularly the Fixed Semantic-Contextual con-1145 figuration, which scores the lowest ViSQOL (3.10) 1146 and PESQ (2.62), despite a relatively high UTMOS 1147 (3.74).1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

Discussion. These results highlight the importance of the temporal alignment strategy in influencing speech reconstruction quality. The superior performance of the **Dynamic-window Contextual** variant demonstrates that token alignment using a dynamic window, where contextual embeddings are adaptively aligned based on token similarity, achieves better semantic grounding and contextual precision.

In contrast, the **Fixed-window** variants suffer from rigid alignment constraints. They fail to capture fine-grained temporal dependencies by enforcing a fixed windowing strategy, which resuls in degraded speech clarity (lower ViSQOL and PESQ). This limitation is especially noticeable in the **Fixed Semantic-Contextual** setup, where the addition of semantic supervision is insufficient to compensate for the strictly aligned contextual embeddings as the fixed window does not account for local content variations.

1169Both Semantic-Contextual variants improve1170UTMOS, indicating that semantic supervision con-1171tributes positively to speech naturalness. However,1172this comes with a trade-off when not paired with1173dynamically aligned contextual guidance, as the1174semantic-only supervision fails to improve content

accuracy.

Overall, these findings underscore that dynamic alignment is essential for effective contextual representation guidance. They also highlight that while semantic supervision enhances fluency and naturalness, it must be combined with flexible alignment mechanisms to avoid compromising content preservation.

C.4 Ablation: Dropout Mask Configuration in Representation Fusion

Setup. We investigate the effect of modality dropout rate on the quality of latent representation fusion. As described in Section 2.3.1, we apply stochastic dropout masks $\mathcal{D}_S, \mathcal{D}_C \in \{0, 1\}^{T' \times D}$ element-wise to the projected semantic (S') and contextual (C') vectors during training:

 $\mathbf{Z}' = \mathbf{Z} + (\mathbf{S}' \odot \mathcal{D}_S) + (\mathbf{C}' \odot \mathcal{D}_C)$ (24)

This stochastic masking prevents FuseCodec from over-reliance on any single modality and encourages the model to learn robust representations.

The selected configuration uses a **10%** dropout rate—i.e., each element in \mathcal{D}_S and \mathcal{D}_C has a 10% chance of being masked to zero during training. We compare this against higher dropout rates: **30%**, **50%**, **70%**, and **90%**.

Results. The best overall performance is achieved with the 10% dropout rate configuration, which achieves the lowest WER (3.99) and WIL (6.45) and the highest STOI (0.95), ViSQOL (3.47), and PESQ (3.13). Increasing the dropout rate to 30–90% leads to the worsening of the most content preservation and speech naturalness metrics. While UTMOS and Similarity remain relatively stable, 50% dropout achieves minor gains in UTMOS (3.66) and Similarity (0.996).

Discussion. These results confirm the impor-
tance of carefully balancing modality dropout12101211

1175 1176 1177

1178 1179 1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

Table 6: Ablation of modality dropout probability during latent representation fusion in FuseCodec. Dropout indicates the stochastic masking rate applied independently to semantic and contextual representations during training. Moderate dropout prevents over-reliance on a single modality, while higher rates degrade multimodal integration. A 10% dropout rate achieves the best trade-off, maximizing intelligibility and perceptual quality.

Model Variant	Dropout	Conte	nt Preser	vation	Speech Naturalness					
	1	WER↓	WIL↓	STOI↑	ViSQOL↑	PESQ↑	UTMOS↑	Similarity↑		
FuseCodec-Fusion	10%	3.99	6.45	0.95	3.47	3.13	<u>3.63</u>	0.995		
FuseCodec-Fusion	30%	4.10	6.63	0.94	3.29	2.96	3.65	0.995		
FuseCodec-Fusion	50%	4.09	6.58	0.94	3.33	2.97	3.66	0.996		
FuseCodec-Fusion	70%	$\overline{4.08}$	6.64	0.93	3.26	2.91	3.63	0.995		
FuseCodec-Fusion	90%	4.15	6.67	<u>0.93</u>	3.26	2.86	3.61	0.995		

during latent fusion and underscore the value of semantic-contextual representation integration. Preserving a sufficient portion of the auxiliary representations by using a small 10% dropout rate achieves the most effective use of semantic and contextual information.

1212

1214

1215

1216

1217

1218

1219

1220

1221

1222

1224

1225

1226

1227

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1243

1244

1245

1246

As the dropout rate increases, the model receives increasingly less additional modality information, reducing its ability to align latent tokens with multimodal supervision. This negatively affects intelligibility (WER, WIL) and perceptual quality (ViSQOL, PESQ).

Interestingly, metrics such as UTMOS and Similarity remain relatively stable or improve at moderate dropout rates (50%), suggesting that prosodic and speaker characteristics are preserved within the base latent representations. However, the loss of some semantic-contextual information comes at the cost of worse content preservation.

Overall, the findings suggest that light dropout (10%) provides the best trade-off, ensuring robust yet expressive multimodal grounding during latent token fusion.

C.5 Ablation: Quntizer Layer Configuration in Semantic-Contextual Guidance

Setup. We study the impact of RVQ layer supervision depth in the distillation objective. Our method, FuseCodec-Distill, uses first-layer supervision, projecting the first-layer RVQ tokens $\mathbf{Q}^{(1)}$ and computing cosine similarity (see Sections 2.3.2 and 2.3.3).

We compare this against an ablated variant, alllayer supervision, which averages the outputs from all eight RVQ layers. We define the averaged RVQ output as:

$$\mathbf{Q}^{(1:8)} = \frac{1}{8} \sum_{i=1}^{8} \mathbf{Q}^{(i)} \in \mathbb{R}^{T' \times D},$$

$$\mathbf{Q}^{(1:8)} = \mathbf{Q}^{(1:8)} \mathbf{W}$$
(25) 1247

In the Global Semantic-Contextual Supervi-1248 sion setting, we apply the all-layer supervision to 1249 the distillation loss as:

$$\mathcal{L}_{\text{distill}} = -\frac{1}{T'} \sum_{t=1}^{T'} \log \sigma \left(\frac{1}{2} \left[\cos \left(\mathbf{Q}_t^{\prime(1:8)}, \tilde{\mathbf{S}}_t \right) + \cos \left(\mathbf{Q}_t^{\prime(1:8)}, \tilde{\mathbf{C}}_t \right) \right] \right)$$
(26)

Similarly, for the Temporally Aligned Contextual Supervision setting, we apply the all-layer supervision to the distillation loss as:

$$\mathcal{L}_{\text{align}} = -\frac{1}{T'} \sum_{t=1}^{T'} \log \sigma \left(\cos \left(\mathbf{Q}_t^{\prime(1:8)}, \mathbf{C}_t^* \right) \right)$$
(27)

Results. Table 7 shows the effect of RVQ su-1256 pervision depth across both distillation configura-1257 tions. For FuseCodec (Distill), which uses Global 1258 Semantic-Contextual Supervision, first-layer super-1259 vision achieves the strongest performance across all 1260 content preservation and naturalness metrics, with 1261 the lowest WER (4.09), WIL (6.60), and highest STOI (0.94), ViSQOL (3.43), PESQ (3.06), UT-1263 MOS (3.65), and Similarity (0.996). Similarly, 1264 FuseCodec (ContexAlign), which uses Temporally 1265 Aligned Contextual Supervision, First-layer su-1266 pervision again achieves stronger results in WER 1267 (4.15), WIL (6.70), ViSQOL (3.18), PESQ (2.85), 1268 and Similarity (0.995). In contrast, using all-layer 1269 supervision leads to consistent degradation across 1270 most metrics in both settings. 1271

1250

1255

1251

1252

1253

Table 7: Ablation of RVQ supervision depth under global (Distill) and temporally aligned (ContexAlign) guidance. **First Layer** indicates supervision is applied only to the first-layer RVQ tokens, while **All Layers** averages representations from all eight RVQ layers before supervision. *Supervising the first-layer RVQ tokens leads to stronger semantic-contextual grounding and improved intelligibility compared to all-layer supervision.*

Model Variant	RVO Laver	Conte	nt Preser	vation	Speech Naturalness				
	itt Q Lujei	WER↓	WIL↓	STOI↑	ViSQOL↑	PESQ↑	UTMOS↑	Similarity↑	
FuseCodec-ContextAlign	First Layer	4.15	6.70	0.93	3.18	2.85	3.65	0.995	
FuseCodec-ContextAlign	All Layers	4.34	7.04	0.93	3.17	2.72	3.65	0.993	
FuseCodec-Distill	First Layer	4.09	6.60	0.94	3.43	3.06	3.65	0.996	
FuseCodec-Distill	All Layers	4.23	6.86	0.93	3.26	2.84	3.61	0.994	

Discussion. The results highlight that the layer at which RVQ tokens are supervised significantly impacts the quality of semantic and contextual guidance during distillation. Supervising the first RVQ layer yields stronger performance, as these tokens encode high-level, abstract representations more aligned with semantic intent and global context. This leads to better linguistic grounding and intelligibility, reflected in improved WER, STOI, and ViSQOL scores.

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1288

1289

1290

1291

1292

1293

1294

1296

1297

1298

1299

1300

1301

1302

1303

1304 1305

1306

1307

1309

In contrast, deeper RVQ layers capture lowerlevel acoustic and residual details, which are less suitable for semantic or contextual alignment. Averaging supervision across all layers matches these fine-grained signals with global ones, impacting the alignment objective. This results in performance drop across content preservation and speech naturalness metrics.

Some naturalness metrics, such as UTMOS and Similarity, remain relatively stable with all-layer supervision, suggesting that speaker identity and prosodic features are distributed throughout the RVQ layers. However, these are insufficient for guiding semantic alignment during distillation.

Overall, applying supervision at the first RVQ layer provides a clearer, more semantically grounded signal, leading to better alignment and overall performance in speech reconstruction.

D Related Work

Recent progress in speech and audio generation has been largely driven by advances in discrete representation learning, neural audio codecs, and language model-based synthesis. VQ-VAE (van den Oord et al., 2018) introduced vector quantization in latent spaces to support symbolic modeling of audio, while HuBERT (Hsu et al., 2021) applied masked prediction over cluster-derived labels to learn contextualized speech features in a self-supervised manner. SoundStream (Zeghidour et al., 2022) proposed a causal adversarially trained codec with residual vector quantization (RVQ) and demonstrated scalable compression at low bitrates. HiFi-Codec (Yang et al., 2023) further improved efficiency by introducing group residual quantization, reducing the number of required codebooks while preserving audio fidelity. On the generative side, AudioLM (Borsos et al., 2023) modeled long-range dependencies in semantic and acoustic tokens using transformer-based language modeling. This approach was extended by VALL-E (Wang et al., 2023), which enabled zero-shot textto-speech synthesis by conditioning on short acoustic prompts and leveraging codec token generation. To improve the suitability of tokenization for language modeling tasks, X-Codec (Ye et al., 2024) integrated semantic embeddings from pretrained models into the quantization pipeline, while LAST (Turetzky and Adi, 2024) learned a tokenizer supervised by a frozen language model to improve downstream ASR and speech generation performance. HiFi-GAN (Kong et al., 2020) introduced multi-period and multi-scale discriminators, enabling high-fidelity waveform synthesis with realtime efficiency.

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

In parallel, codec designs have evolved to im-1336 prove training stability and perceptual quality. En-1337 Codec (Défossez et al., 2022) introduced a GAN-1338 based codec architecture with multi-loss balancing 1339 and spectrogram-based discrimination, setting a 1340 new benchmark for real-time low-bitrate synthe-1341 sis. BigCodec (Xin et al., 2024) scaled the VQ-1342 VAE framework and showed that a single large 1343 codebook could achieve near-human perceptual 1344 quality at 1 kbps. DAC (Kumar et al., 2023) pro-1345 posed refinements to residual quantization, such 1346 as factorized and normalized codebooks, and in-1347 troduced advanced discriminators to improve qual-1348

1349	ity under bitrate constraints. More recent work
1350	has focused on improving token expressiveness for
1351	downstream tasks. SpeechTokenizer (Zhang et al.,
1352	2024) demonstrated that hierarchical quantization
1353	improves resynthesis and zero-shot TTS, while
1354	DM-Codec (Ahasan et al., 2024) aligned quantiza-
1355	tion layers with pretrained speech and text models
1356	to reduce WER and enhance contextual fidelity. Fi-
1357	nally, NaturalSpeech 3 (Ju et al., 2024) introduced
1358	a factorized codec for disentangling prosodic and
1359	acoustic attributes in speech, and Moshi (Défossez
1360	et al., 2024) unified ASR and TTS in a streaming,
1361	full-duplex transformer model operating on jointly
1362	learned semantic and acoustic tokens.

1363 E Qualitative Comparison



Sample 1

(click image to play audio)

Sample 2

Figure 2: Qualitative speech reconstruction results comparing our method to multiple baselines. Each cell shows the spectrogram output for two samples; clicking an image plays the corresponding audio.