

000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 CHANGANYA NDIMI: CODE-SWITCHED SPEECH GENERATION WITH A DIFFUSION PRIOR AND LINGUISTIC CONSTRAINTS

Anonymous authors

Paper under double-blind review

ABSTRACT

We address code-switched speech generation by *editing* monolingual utterances with a pretrained diffusion-based speech model guided by linguistic constraints. Our method requires no parallel code-switched data. Instead, generation is conditioned on two differentiable modules—a multilingual language classifier and a contrastively trained segment encoder—that jointly guide where to insert semantically coherent, sociolinguistically appropriate foreign segments. During reverse diffusion, the system iteratively refines noisy speech representations, performing targeted segment substitutions while preserving fluency, prosody, and speaker identity.

On a semantically aligned corpus spanning five African languages from three language families, our approach achieves strong performance: segment-level COMET 0.815, LaBSE similarity 0.880, and 6.7% Equal Error Rate (EER) for speaker identity preservation. The model also reproduces natural code-switching patterns—frequency, temporal distribution, and alternation rates—without explicit supervision for such behaviors. To our knowledge, this is the first system to enable *controlled multilingual infusion within a single utterance*, highlighting guided diffusion as a flexible, plug-and-play framework for low-resource multilingual speech generation. Audio samples are available at: <https://github.com/codeSwitchLugha/CodeSwitch>.

1 INTRODUCTION

Code-switching—the fluid alternation between languages within an utterance—is widespread among African speakers Biswas et al. (2022); Sitaram et al. (2019). Yet most speech technologies (ASR, S2ST, SLU, speech LLMs) still rely mainly on monolingual data, as high-quality code-switched corpora are scarce. Code-switched speech synthesis remains relatively underexplored, and collecting spontaneous code-switched audio is costly and often yields unnatural speech Tarunesh et al. (2021); Hsu et al. (2023). We address this gap with a method that transforms monolingual corpora into realistic code-switched utterances via minimal, semantically coherent edits. Our approach samples from a constrained denoising diffusion model (DDPM): starting from noise, the model iteratively refines an utterance while two differentiable controllers guide generation:

- (a) $c_1(x, y)$: a language-ID (LID) controller that decides *where* and *how much* to switch;
- (b) $c_2(x, y)$: a multilingual segment encoder that determines *what* foreign-language content to insert by swapping in semantically matched segments.

We follow the *plug-and-play diffusion* paradigm: instead of retraining the generative model, we keep the diffusion prior $p(x)$ frozen and attach external constraint modules that steer sampling at test time. Formally,

$$p(x | y) \propto p(x) C(x, y), \quad C(x, y) = c_1(x, y) c_2(x, y), \quad (1)$$

where x denotes an utterance-level waveform and y encodes the infusion specification (host/foreign language set, switch prior, optional source semantics, retrieval index). The prior $p(x)$ models natural speech (speaker identity, prosody), while $C(x, y)$ modulates *when/how much* to leave the host language (c_1) and *what* foreign content to insert (c_2). Guidance is implemented as a short, time-ramped penalty during sampling; the full derivation is given in App. B.

054 Experiments on four African languages plus English show that the synthesized speech is fluent, se-
 055 mantically aligned, and speaker-consistent according to automatic metrics (SacreBLEU, BERTScore,
 056 COMET, LaBSE, ECAPA-TDNN EER) and human judgments. The method can supply code-
 057 switched data for low-resource speech recognition, speech-to-speech translation, and multilingual
 058 LLM training. Our contributions are:

059

- 060 **1. Retrieval-augmented, plug-and-play diffusion** for code-switched speech that requires no parallel
 061 CS data and offers controllable *where/what* switching via external controllers on a frozen diffusion
 062 prior.
- 063 **2. Two complementary controllers:** LID-based switching (c_1) and retrieval-based semantic infusion
 064 (c_2) with late-commit and blend-and-write-back schedules.
- 065 **3. Evaluation suite** covering semantic fidelity, speaker consistency, prosody continuity at switch
 066 boundaries, code-switch structure, and human ratings.
- 067 **4. Empirical results on five languages** (Swahili, Luo, Kikuyu, Nandi, English) showing fluent,
 068 semantically aligned, speaker-consistent code-switching and downstream utility.

070 2 METHOD

072 2.1 PROBLEM FORMULATION

074 Our goal is to sample code-switched utterances from the constrained posterior in Eq. 1. Here, x
 075 denotes a waveform and y is an *infusion specification* that encodes how code-switching should occur.
 076 We write $y = (\mathcal{S}_{\text{inf}}, \pi_{\text{switch}}, \phi_{\text{src}}, \mathcal{R})$, where \mathcal{S}_{inf} is the allowed foreign-language set, π_{switch} is a
 077 prior over *where/how much* to switch, ϕ_{src} is a semantic representation of the host utterance, and \mathcal{R}
 078 is a retrieval index over foreign-language segments. The constraint term $C(x, y)$ (Eq. 1) factorizes as
 079 $c_1(x, y)c_2(x, y)$, where c_1 encourages plausible switch locations given π_{switch} and the LID model,
 080 and c_2 enforces semantic consistency between the edited utterance and the host. The prior $p(x)$
 081 is implemented by a pretrained DDPM vocoder (Sec. 2.2), which ensures natural speaker identity,
 082 prosody, and acoustic realism.

083 Direct inference in $p(x | y)$ is intractable because the diffusion prior introduces a sequence of latent
 084 noise variables. We therefore augment x with diffusion latents $h = \{x_1, \dots, x_T\}$ and work with the
 085 joint model $p(x, h) = p(h)p(x | h)$, leading to the variational free-energy objective

$$086 F(q) = \text{KL}(q(x, h) \| p(x, h)) - \mathbb{E}_{q(x)}[\log C(x, y)], \quad (2)$$

088 where $q(x, h)$ is a variational distribution over utterances and diffusion trajectories. The KL term
 089 favors utterances that are likely under the diffusion prior $p(x, h)$, while $-\mathbb{E}_{q(x)}[\log C(x, y)]$ steers
 090 samples toward satisfying the code-switching constraints. A step-by-step derivation of Eq. 2 from
 091 Eq. 1 and the DDPM joint model $p(x, h)$ is given in Appendix A.

092 2.2 FREE-ENERGY OBJECTIVE IN THE DDPM FRAMEWORK

094 We use a standard denoising diffusion probabilistic model (DDPM) Ho et al. (2020) as a frozen
 095 speech prior. The forward process gradually corrupts clean speech x_0 to noise, $q(x_t | x_0) =$
 096 $\mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I)$, $t = 1, \dots, T$, and the reverse process is parameterized as $p_\theta(x_{t-1} |$
 097 $x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$, $x_T \sim \mathcal{N}(0, I)$. As in Ho et al. (2020), maximizing $\log p_\theta(x_0)$
 098 is equivalent (up to constants) to minimizing the noise-prediction loss

$$100 \mathcal{L}_{\text{DDPM}} = \mathbb{E}_{t, \epsilon_0} \left[\|\hat{e}_\theta(x_t, t, M) - \epsilon_0\|_2^2 \right], \quad x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_0, \quad (3)$$

102 where $\epsilon_0 \sim \mathcal{N}(0, I)$ and M denotes conditioning inputs (speaker, text, etc.).

103 To incorporate linguistic constraints without retraining the prior, we target the constrained posterior
 104 in Eq. 1, where $p(x)$ is the DDPM prior, y encodes the infusion specification, and c_1, c_2 are our
 105 controllers (Sec. 1). Direct inference in Eq. 1 is intractable because $p(x)$ is defined via the latent
 106 diffusion trajectory $h = \{x_1, \dots, x_T\}$. Following Chung et al. (2023), we adopt a mode-seeking
 107 variational family $q(x) = \delta(x - \eta)$, so that the free-energy objective in Eq. 2 reduces to a function of
 a single clean sample η .

108 Reparameterizing $x_t = \sqrt{\bar{\alpha}_t} \eta + \sqrt{1 - \bar{\alpha}_t} \epsilon_0$ and sampling $t \sim U(1, T)$, we obtain the practical
 109 plug-and-play objective used in our guided sampler:
 110

$$111 \quad F(\eta) = \mathbb{E}_{t, \epsilon_0} \left[\|\hat{\epsilon}_\theta(x_t, t, M) - \epsilon_0\|_2^2 \right] - \log C(\eta, y), \quad (4)$$

113 where the first term pushes η toward high-likelihood speech under the prior and $-\log C(\eta, y)$ adds
 114 soft guidance from the switch and semantic controllers. In practice, this corresponds to standard
 115 DDPM reverse updates augmented with constraint gradients, realizing plug-and-play code-switching:
 116 the diffusion prior $p(x)$ remains frozen and the linguistic controllers steer generation at test time.
 117 Appendix B provides the complete derivation from the joint model $p(x, h)$ to Eq. 4.
 118

119 3 DIFFUSION-BASED CODE-SWITCHING MODEL (DCSM)

120 3.1 CONSTRAINT $c_1(x, y)$: LANGUAGE IDENTIFICATION AND INFUSION DECISION

123 The first constraint $c_1(x, y)$ decides whether a *monolingual* host utterance x_j should undergo foreign-
 124 language infusion, and if so, how strongly. Utterances that already contain substantial foreign
 125 material are left as they are, while mostly clean monolingual speech is treated as a candidate for
 126 code-switching.

127 **Segment-level foreignness.** We segment x_j into n short spans $\{s_{x_j}^{(i)}\}_{i=1}^n$ (e.g., log-Mel windows).
 128 Each span is fed to a frozen multilingual LID classifier f_{cl} , which outputs a posterior over the
 129 host language ℓ_{mono} and the infusion-eligible languages $\mathcal{S}_{\text{inf}}(y) = \{\ell_1, \dots, \ell_m\}$: $p(\ell \mid s_{x_j}^{(i)}) =$
 130 $f_{\text{cl}}(s_{x_j}^{(i)})_\ell, \quad \ell \in \mathcal{S}_{\text{all}} = \{\ell_{\text{mono}}\} \cup \mathcal{S}_{\text{inf}}(y)$. For each segment we define a *foreignness score*
 131 $u^{(i)} = \sum_{\ell \in \mathcal{S}_{\text{inf}}(y)} p(\ell \mid s_{x_j}^{(i)}) \in [0, 1]$, and aggregate these into an utterance-level measure
 132

$$134 \quad P_{\text{foreign}}(x_j) = \frac{1}{n} \sum_{i=1}^n u^{(i)} \in [0, 1]. \quad (5)$$

137 Here, $P_{\text{foreign}}(x_j) \approx 0$ indicates that the classifier sees x_j as strongly monolingual in ℓ_{mono} , while
 138 larger values signal that many segments already exhibit foreign-language characteristics (existing
 139 code-switches, borrowings, or noisy labels).

140 **Global controller and gating.** We use $P_{\text{foreign}}(x_j)$ to control whether and how much the utterance
 141 should be edited. Given a target switch rate $\pi_{\text{switch}} \in [0, 1]$, we define $c_1(x_j, y) = \sigma(\alpha_1 [\pi_{\text{switch}} -$
 142 $P_{\text{foreign}}(x_j)])$, with sharpness $\alpha_1 > 0$ and sigmoid $\sigma(\cdot)$. When the utterance is *less* foreign than
 143 the target ($P_{\text{foreign}}(x_j) < \pi_{\text{switch}}$), $c_1(x_j, y) \approx 1$ and the constraint encourages infusion; when it
 144 is already *more* foreign than desired, $c_1(x_j, y)$ shrinks toward 0 and suppresses further switching.
 145 For stricter control, we optionally use a hard gate $z_{\text{infuse}} = \mathbb{1}(P_{\text{foreign}}(x_j) \leq \tau)$, with tolerance
 146 $\tau \in [0, 1]$. Only utterances whose foreignness is below τ are considered for code-switching; those
 147 with $P_{\text{foreign}}(x_j) > \tau$ are left untouched.

149 **Local soft masks for segment selection.** The same foreignness scores also provide soft, per-
 150 segment priorities for where infusion should happen. We define local gates $g^{(i)} = \sigma(\alpha_1 [\pi_{\text{switch}} -$
 151 $u^{(i)}])$, so segments that look *more* monolingual (low $u^{(i)}$) receive higher weights $g^{(i)}$ and are
 152 preferred as candidates for replacement inside the infusion constraint c_2 . At inference time, one
 153 may derive hard gates $z^{(i)} \in \{0, 1\}$ from $g^{(i)}$, while keeping gradients through the sigmoid during
 154 backprop.

156 **Contribution to the guided objective.** In the guided sampling objective, c_1 contributes via the
 157 scalar penalty
 158

$$159 \quad \mathcal{L}_{c_1}(x_j, y) = -\log c_1(x_j, y) = -\log \sigma(\alpha_1 [\pi_{\text{switch}} - P_{\text{foreign}}(x_j)]), \quad (6)$$

160 which is added to the guided objective in Eq. 4 with weight λ_1 . This directly ties the LID-based
 161 controller to the diffusion trajectory: utterances that are too monolingual relative to π_{switch} are

162 nudged toward more foreign infusion, while utterances that are already highly foreign are protected
 163 from further editing. The classifier f_{cl} is pretrained and frozen during diffusion. It is trained to detect
 164 segment-level foreign-language presence using a standard cross-entropy objective; we provide the
 165 full training loss and label specification in Appendix H.
 166

167 **3.2 CONSTRAINT $c_2(x, y)$: FOREIGN SEGMENT INFUSION**
 168

169 At each denoising step, we edit a *single* host span chosen by c_1 . Let $i^* = \arg \max_i g^{(i)}$ be the
 170 selected segment in utterance x_j ; constraint c_2 replaces $s_{x_j}^{(i^*)}$ with a foreign segment that is (1)
 171 semantically similar and (2) prosodically compatible.
 172

173 **Semantic retrieval.** We encode the source segment with a frozen multilingual encoder, $q =$
 174 $f_{enc}(s_{x_j}^{(i^*)})$, and query a FAISS-based *vector database* Johnson et al. (2019) \mathcal{D} . The database
 175 stores pre-segmented foreign-language spans, each with an ℓ_2 -normalized embedding $f_{enc}(s_{y_k}^{(m)})$ and
 176 metadata (duration, onset, language tag). During retrieval, we restrict candidates to the infusion-
 177 eligible language set $\mathcal{S}_{\text{inf}}(y)$ and select the best match under cosine similarity:
 178

$$179 \quad m^* = \arg \max_{\substack{s_{y_k}^{(m)} \in \mathcal{D} \\ \ell(s_{y_k}^{(m)}) \in \mathcal{S}_{\text{inf}}(y)}} \text{Sim}(q, f_{enc}(s_{y_k}^{(m)})), \quad s^* = s_{y_k}^{(m^*)}.$$

$$180$$

$$181$$

$$182$$

183 Appendix C describes an optional soft top- M variant with a temperature schedule that anneals from a
 184 mixture to hard top-1 retrieval.
 185

186 **Prosodic compatibility.** To avoid audible glitches, we require the retrieved segment to match the
 187 host in duration and timing. We estimate an expected candidate duration \hat{d} from global speech-rate
 188 and prosodic statistics of the host and candidate utterances, then define hinge penalties for duration
 189 and onset mismatch:
 190

$$191 \quad \mathcal{L}_{\text{dur}} = \max\left(0, \frac{|d_{y_k}^{(m^*)} - \hat{d}| - \lambda_d \hat{d}}{\hat{d}}\right), \quad \mathcal{L}_{\text{on}} = \max\left(0, \frac{|\mathcal{O}_{y_k}^{(m^*)} - \mathcal{O}_{x_j}^{(i^*)}| - \Delta\tau}{\Delta\tau}\right),$$

$$192$$

$$193$$

194 where $\lambda_d \in [0, 1)$ and $\Delta\tau$ are tolerance parameters. Full expressions for \hat{d} , λ_d , and $\Delta\tau$ are given in
 195 Appendix D.
 196

197 **Semantic and contextual consistency.** Beyond segment-level similarity, we require the injected
 198 segment to be coherent with its local context. We define a semantic loss
 199

$$200 \quad \mathcal{L}_{\text{sem}} = -\text{Sim}\left(f_{enc}(s_{x_j}^{(i^*)}), f_{enc}(s^*)\right),$$

$$201$$

202 and a contextual loss over neighboring host segments $\mathcal{N}(i^*)$:
 203

$$204 \quad \mathcal{L}_{\text{ctx}} = -\frac{1}{|\mathcal{N}(i^*)|} \sum_{s' \in \mathcal{N}(i^*)} \text{Sim}\left(f_{enc}(s^*), f_{enc}(s')\right).$$

$$205$$

206 The neighborhood definition and window size are specified in Appendix E.
 207

208 **Per-step loss and blend-and-write-back.** The infusion loss at this step is
 209

$$210 \quad \mathcal{L}_{c_2} = g^{(i^*)} [\alpha_{\text{sem}} \mathcal{L}_{\text{sem}} + \alpha_{\text{ctx}} \mathcal{L}_{\text{ctx}} + \alpha_{\text{pro}} (\mathcal{L}_{\text{dur}} + \mathcal{L}_{\text{on}})],$$

$$211$$

212 with fixed weights $\alpha_{\text{sem}} = \alpha_{\text{ctx}} = 1$ and $\alpha_{\text{pro}} = 0.1$. We then apply a time-dependent blend-and-write-
 213 back update $s_{x_j}^{(i^*)} \leftarrow (1 - \rho_t) s_{x_j}^{(i^*)} + \rho_t s^*$, where the ramp ρ_t increases over reverse-diffusion steps
 214 so that semantic edits are committed only after coarse acoustic structure has stabilized. Gradients
 215 flow back to the clean sample η via the inputs to f_{enc} , while both f_{enc} and the retrieval index \mathcal{D} remain
 frozen.
 216

216 3.3 GUIDED SAMPLING OBJECTIVE AND INFERENCE
217

218 During generation, we implement the free-energy objective in Eq. equation 4 by optimizing the
219 clean sample η while keeping the DDPM vocoder and constraint networks fixed. At a given reverse-
220 diffusion step t , the guided loss is

$$221 \quad \mathcal{L}_{\text{step}}(t) = \|\hat{\epsilon}_{\theta}(x_t, t, M) - \epsilon_0\|_2^2 + \lambda_1 \mathcal{L}_{c_1} + \lambda_2(t) \mathcal{L}_{c_2}, \quad (7)$$

223 where $x_t = \sqrt{\bar{\alpha}_t} \eta + \sqrt{1 - \bar{\alpha}_t} \epsilon_0$ and $\epsilon_0 \sim \mathcal{N}(0, I)$ as in Eq. 3. The first term is the DDPM noise-
224 prediction loss; the constraint losses \mathcal{L}_{c_1} and \mathcal{L}_{c_2} are defined in Sections 3.1–3.2, and $\lambda_2(t) \in [0, 1]$ is
225 a time-dependent weight that ramps constraint guidance over diffusion steps. The DDPM parameters
226 θ remain frozen; gradients flow only to η .

227 To regulate when foreign segments are infused, we use two schedules: (i) a blending coefficient
228 ρ_t that increases over denoising steps, so that semantic edits are only fully committed once coarse
229 acoustic structure has stabilized; and (ii) the guidance weight $\lambda_2(t)$, which delays the effect of \mathcal{L}_{c_2}
230 until early noise has largely dissipated. These mechanisms stabilize generation and support gradual
231 linguistic modulation.

232 Inference proceeds by iteratively denoising a Gaussian sample while selectively modifying one span
233 per step: at each t , c_1 selects a candidate segment, c_2 retrieves and blends in a foreign segment using
234 ρ_t , and we update η by taking a gradient step on $\mathcal{L}_{\text{step}}(t)$. Over time, this process converges to a
235 fluent, code-switched utterance. The full inference procedure is given in Appendix F.

236 3.4 SPEAKER IDENTITY HARMONIZATION
237

239 To standardize timbre, pitch, and rhythm across edited spans *while preserving content*, we apply a
240 short identity-harmonization pass with the *pretrained, frozen* DDPM prior. Given a monolingual
241 reference utterance x_{mono} and the current code-switched sample x , we extract speaker/prosody
242 descriptors

$$243 \quad \phi_{\text{spk}} = \text{ECAPA}(x_{\text{mono}}), \quad \phi_{\text{mel}} = \text{MelStats}(x),$$

244 where ECAPA(\cdot) is a frozen ECAPA-TDNN and MelStats(\cdot) denotes summary statistics (e.g., mean
245 and variance) of log-Mel features. We then run a shallow denoising refinement

$$246 \quad x_{\text{final}} = \text{Refine}_{\theta}(x \mid \phi_{\text{spk}}, \phi_{\text{mel}}), \quad (8)$$

248 using $T_{\text{ref}} = 150$ diffusion steps with a low-noise schedule (late timesteps only).

249 During this refinement we disable segment edits by setting $\lambda_2(t) = 0$ and keep the LID-based rate
250 term weak (small λ_1), so as not to alter semantics or the code-switch pattern. In practice, conditioning
251 can be implemented via feature concatenation or FiLM-style modulation inside the DDPM U-Net,
252 with θ kept fixed. Alternative variants (e.g., adding a cosine speaker-embedding penalty \mathcal{L}_{id} computed
253 by a frozen ECAPA on sliding windows) yield similar improvements; see Appendix G for details.

254 4 EVALUATION
255256 4.1 DATASET
257

259 We collected a proprietary speech dataset from the Kenya Broadcasting Corporation (KBC), which
260 operates 11 radio stations delivering aligned news content across multiple Kenyan languages. News
261 bulletins are authored in English and translated into local languages, then read aloud by native
262 speakers—yielding semantically aligned monolingual utterances across languages. Our corpus spans
263 2018–2023 and focuses on the 7 p.m. bulletins, which are typically the most content-rich. We
264 retain advertisements, presenter introductions, and other ambient segments to preserve real-world
265 variability. Each bulletin is segmented using an over-segmentation VAD pipeline Duquenne et al.
266 (2021), producing speech units bounded by silence and ranging from 3 to 20 seconds.

267 We focus on five languages—Swahili, Luo, Kikuyu, Nandi, and English—selected for their regional
268 and typological diversity. Swahili and Kikuyu belong to the Niger–Congo phylum, while Luo and
269 Nandi are Nilo–Saharan; English, though non-indigenous, functions as a lingua franca and appears
in every station’s programming. Table 1 summarizes dataset statistics by language, and Table 2

270 reports segment-level details. Bulletins are read by multiple presenters, capturing a range of accents,
 271 pitch ranges, and prosodic patterns, which improves generalization for synthesis, translation, and
 272 recognition models. We apply a 70/30 split of segments per language for training and evaluation, and
 273 zero-pad each segment to a fixed length of 20 seconds for model training.

274 In addition to news bulletins, we collected 7,255 naturally occurring code-switched Swahili–Luo
 275 utterances from radio call-in segments. For evaluation, each utterance was manually re-recorded in a
 276 monolingual version: 4,044 were rendered in Swahili and 3,211 in Luo, depending on the dominant
 277 language of the original speaker.

280 **Table 1: Monolingual speech dataset sum-
 281 mary.**

Language	Family	Daily Bulletins (2018–2023)	Total Hours
Swahili	Niger-Congo	2190	1353
Luo	Nilo-Saharan	2190	1284
Kikuyu	Niger-Congo	2190	1304
Nandi	Nilo-Saharan	2190	1256
English	–	2190	1206

282 **Table 2: Segment statistics per language.**

Language	Avg. Segment Length (s)	Total Segments
Luo	16.5	601,112
Nandi	17.1	594,503
Kikuyu	15.6	643,001
English	15.0	665,578
Swahili	15.2	638,944

288 4.2 TOOLS AND RESOURCES

289 **Tools and Models.** Our system relies on four frozen modules: (i) a multilingual segment-level LID
 290 classifier; (ii) a contrastively trained speech encoder for semantic retrieval and contextual matching;
 291 (iii) the SegUniDiff model for conditional speech generation and refinement; and (iv) lightweight
 292 ASR and MT systems used only for evaluation. Full architectural specifications, training details, and
 293 performance metrics for all components are provided in Appendix H.

296 4.3 EVALUATION METRICS AND FRAMEWORK

297 We evaluate both segment-level and utterance-level code-switched speech using four complementary
 298 metrics that capture lexical accuracy, semantic preservation, and cross-lingual coherence: Sacre-
 299 BLEU Post (2018) for surface correspondence, BERTScore Zhang et al. (2019) for contextual
 300 similarity, COMET Rei et al. (2020) for semantic adequacy, and LaBSE cosine similarity Feng et al.
 301 (2022) for cross-lingual embedding alignment. For each metric, we report mean scores with 95%
 302 confidence intervals obtained through bootstrap resampling. Full metric definitions, evaluation setups,
 303 and the resampling protocol are provided in Appendix I.

305 4.4 SEGMENT-LEVEL EVALUATION WITH CONFIDENCE INTERVALS

306 We assess segment-level semantic fidelity using the metrics defined in Section 4.3. Our evaluation
 307 covers 8,500 synthetic utterances and 7,255 naturally occurring code-switched utterances (primarily
 308 Swahili–Luo) collected from radio call-ins. For each utterance x , we extract VAD-based segments
 309 (Section 4.1), apply the LID classifier f_{cl} to detect foreign segments, and pair each detected code-
 310 switched segment $s_{x_c}^{(k)}$ (segment k of the code-switched utterance x_c) with its corresponding source
 311 segment $s_x^{(k)}$. Segment pairs are transcribed with ASR, translated to English, and scored on semantic
 312 fidelity. Table 3 reports results with 95% confidence intervals from 1,000 bootstrap samples; the
 313 resampling protocol is detailed in Appendix J.

314 **Table 3: Segment-level evaluation of synthetic and natural code-switched utterances across four
 315 metrics. Scores include 95% confidence intervals from 1,000 bootstrap samples.**

Source	SacreBLEU (↑)	BERTScore (↑)	COMET (↑)	LaBSE (↑)
Swahili (Synthetic)	38.4 [36.9, 39.6]	0.814 [0.809, 0.818]	0.831 [0.822, 0.843]	0.890 [0.882, 0.897]
Luo (Synthetic)	35.7 [34.2, 37.1]	0.805 [0.801, 0.810]	0.809 [0.798, 0.819]	0.876 [0.869, 0.884]
Kikuyu (Synthetic)	36.8 [35.0, 38.4]	0.808 [0.802, 0.813]	0.817 [0.806, 0.828]	0.881 [0.874, 0.888]
Nandi (Synthetic)	34.5 [33.1, 35.7]	0.801 [0.795, 0.806]	0.804 [0.794, 0.814]	0.872 [0.865, 0.880]
Luo (Natural)	36.2 [35.4, 37.0]	0.822 [0.818, 0.841]	0.833 [0.827, 0.843]	0.885 [0.879, 0.891]
Swahili (Natural)	39.0 [38.2, 39.7]	0.836 [0.831, 0.841]	0.845 [0.841, 0.864]	0.898 [0.892, 0.904]
Average (Synthetic)	36.4 [35.2, 37.4]	0.807 [0.804, 0.810]	0.815 [0.807, 0.823]	0.880 [0.873, 0.886]

Synthetic utterances closely approximate natural ones across all metrics. For Swahili and Luo, the COMET gaps between synthetic and natural segments are only 0.014 and 0.024, respectively; LaBSE gaps are similarly small (0.008 and 0.009). These results show that our guided segment substitution preserves cross-lingual semantics without requiring naturally code-switched training data. The narrow confidence intervals further indicate that performance is stable and not driven by outliers.

4.5 UTTERANCE-LEVEL EVALUATION WITH MASKED AND FULL VARIANTS

To assess overall fluency, inter-segment coherence, and disruptions introduced by code-switching, we evaluate at the utterance level. This complements segment-level analysis by capturing prosodic mismatches, semantic drift, and syntactic incongruities that emerge only in longer contexts. We evaluate both **synthetic** code-switched utterances and **natural** ones collected from Swahili and Luo radio call-ins. For each code-switched utterance x_c , we apply the language classifier f_{cl} to identify foreign-language spans and evaluate under two variants:

- **Full reconstruction (unmasked):** Foreign segments are translated into the source language and reinserted, yielding a reconstructed monolingual utterance x_r .
- **Masked evaluation:** Foreign segments are removed from x_c , yielding x_m , which isolates preservation of the monolingual portions.

Each variant is compared to the clean reference utterance using ASR+MT to obtain transcriptions, followed by SacreBLEU, BERTScore, COMET, and LaBSE similarity. Table 4 reports mean scores with 95% confidence intervals; the resampling protocol is detailed in Appendix K.

Table 4: Utterance-level evaluation of code-switched speech. Each source shows masked and full scores with 95% confidence intervals.

Source	Type	SacreBLEU (\uparrow)	BERTScore (\uparrow)	COMET (\uparrow)	LaBSE (\uparrow)
Swahili (Synthetic)	Full	36.6 [35.3, 37.8]	0.762 [0.757, 0.766]	0.669 [0.660, 0.681]	0.882 [0.875, 0.888]
	Masked	34.9 [33.5, 36.0]	0.737 [0.732, 0.737]	0.642 [0.631, 0.655]	0.854 [0.846, 0.860]
Luo (Synthetic)	Full	33.9 [32.6, 35.3]	0.753 [0.747, 0.758]	0.647 [0.637, 0.657]	0.871 [0.864, 0.878]
	Masked	32.2 [30.9, 33.6]	0.728 [0.722, 0.733]	0.620 [0.609, 0.631]	0.844 [0.837, 0.854]
Kikuyu (Synthetic)	Full	35.0 [33.4, 36.4]	0.756 [0.750, 0.761]	0.655 [0.644, 0.666]	0.876 [0.870, 0.882]
	Masked	33.3 [31.8, 34.6]	0.731 [0.725, 0.736]	0.628 [0.617, 0.639]	0.850 [0.843, 0.857]
Nandi (Synthetic)	Full	32.7 [31.5, 33.9]	0.749 [0.743, 0.756]	0.643 [0.633, 0.653]	0.869 [0.861, 0.875]
	Masked	31.2 [29.9, 32.4]	0.724 [0.717, 0.731]	0.615 [0.604, 0.625]	0.841 [0.834, 0.849]
Swahili (Natural)	Full	37.3 [36.5, 38.1]	0.785 [0.780, 0.791]	0.701 [0.692, 0.710]	0.896 [0.890, 0.902]
	Masked	35.5 [34.7, 36.3]	0.760 [0.755, 0.765]	0.667 [0.658, 0.676]	0.867 [0.860, 0.873]
Luo (Natural)	Full	35.2 [34.1, 37.3]	0.772 [0.767, 0.778]	0.682 [0.673, 0.691]	0.884 [0.878, 0.891]
	Masked	33.6 [32.5, 35.6]	0.745 [0.740, 0.751]	0.654 [0.645, 0.662]	0.856 [0.850, 0.862]
Average (Synthetic)	Full	34.5 [33.5, 35.6]	0.755 [0.751, 0.759]	0.653 [0.645, 0.661]	0.874 [0.870, 0.878]
	Masked	32.9 [31.9, 34.0]	0.730 [0.726, 0.734]	0.626 [0.617, 0.635]	0.847 [0.843, 0.852]

As expected, utterance-level scores are lower than segment-level ones (Table 3), since longer contexts expose more opportunities for prosodic and discourse mismatches. Nevertheless, the model retains strong fluency and coherence: synthetic full-utterance scores are close to their natural counterparts (e.g., Swahili COMET 0.669 vs. 0.701 and LaBSE 0.882 vs. 0.896). The masked variant shows that unaltered monolingual content is largely preserved, with only modest drops relative to the full reconstruction. Taken together, these results indicate that our guided diffusion model can introduce foreign segments while maintaining global utterance quality without access to naturally code-switched training data.

5 SPEAKER IDENTITY VERIFICATION

To evaluate whether code-switched speech maintains consistent speaker identity, we use an automatic verification framework based on ECAPA-TDNN embeddings Desplanques et al. (2020). The model is trained on 732 speakers and produces fixed-length embeddings from short segments. For each generated code-switched utterance x_c , we apply VAD segmentation (Section 4.1) and extract speaker embeddings $f(s_{x_c}^{(i)})$ for each segment $s_{x_c}^{(i)}$.

We then compute cosine similarity between all intra-utterance segment pairs ($i \neq j$), treating these as **genuine pairs** that should correspond to a single speaker. **Impostor pairs** are formed by pairing segments from different utterances (i.e., $s_{x_c}^{(i)}$ and $s_{x'_c}^{(j)}$), assuming different speaker prompts.

378 Speaker consistency is quantified using:
 379

- 380 • **Average cosine similarity** (\uparrow) between genuine and impostor pairs.
- 381 • **Equal Error Rate (EER)** (\downarrow): the point where false accept and false reject rates intersect.

383 **Table 5: Speaker verification results for code-switched utterances.**

384 Source	385 Avg. Cosine Similarity (Genuine) \uparrow	386 Avg. Cosine Similarity (Impostor) \downarrow	387 Equal Error Rate (EER %) \downarrow
388 Swahili (Synthetic)	0.872	0.432	6.5
389 Luo (Synthetic)	0.861	0.418	7.2
Kikuyu (Synthetic)	0.868	0.427	6.8
Nandi (Synthetic)	0.854	0.411	7.6
Swahili (Natural)	0.903	0.391	3.6
Luo (Natural)	0.870	0.430	5.1
Average	0.868	0.426	6.7

390 Table 5 shows that speaker identity is generally preserved across code-switched utterances. As
 391 expected, **natural** utterances perform best, with higher genuine-pair similarity and lower EERs
 392 (e.g., Swahili: 0.903 similarity, 3.6% EER). Synthetic utterances also score well, with EERs in the
 393 6.5–7.6% range and genuine similarities above 0.85. The relatively small gap between synthetic and
 394 natural conditions suggests that our model retains speaker traits across substituted segments, enabling
 395 code-switching that is both semantically faithful and vocally consistent.
 396

397 398 5.1 CODE-SWITCHING PATTERNS ACROSS SOURCE LANGUAGES

400 We analyze generated utterances along four structural dimensions: (i) switching frequency, (ii)
 401 distribution of inserted languages, (iii) temporal position of switches, and (iv) alternation points
 402 between languages.

403 **Switching frequency.** We sample 2,000 synthetic utterances per source language, segment them
 404 via VAD (Section 4.1), and label each segment with the pretrained classifier f_{cl} . We then compute the
 405 average proportion of foreign segments per utterance and compare to natural call-in data (Table 6).

408 **Table 6: Average percentage of foreign segments per utterance (higher = more frequent code-
 409 switching).**

410 Source Language	411 Foreign Segment Rate
412 Swahili (Synthetic)	4.8%
413 Luo (Synthetic)	4.2%
Kikuyu (Synthetic)	4.4%
Nandi (Synthetic)	3.9%
Swahili (Natural)	5.3%
Luo (Natural)	3.2%
Average (Synthetic)	4.3%

415 Synthetic utterances exhibit realistic switching rates, typically within 1–1.5 percentage points of natu-
 416 ral baselines. Swahili shows the highest frequency in both synthetic and natural settings, consistent
 417 with its role as a lingua franca.

419 **Inserted language and temporal position.** Table 7 shows normalized insertion frequencies by
 420 source language. For Swahili (which allows all other languages as infusion targets), insertions are
 421 diverse. For restricted sources (Luo, Kikuyu, Nandi), Swahili is preferred over English, reflecting
 422 both phonological compatibility and its empirical prevalence in the data.

424 Table 8 reports the proportion of foreign segments per utterance quarter (Q1: start, Q4: end). Synthetic
 425 patterns track natural ones: Swahili places more foreign material toward the end of the utterance,
 426 whereas Luo shows a flatter distribution.

428 **Alternation points.** We define the alternation rate as the proportion of segment boundaries where
 429 the language label changes:

431
$$\text{Alternation Rate} = \frac{\#\{\text{boundaries where language changes}\}}{\#\{\text{segment boundaries}\}}.$$

432

433 Table 7: Distribution (%) of inserted languages
434 per source utterance.

Insert \ Source	Swahili	Luo	Kikuyu	Nandi
Insert	Swahili	Luo	Kikuyu	Nandi
English	21.2%	46.7%	48.9%	44.5%
Luo	28.5%	—	—	—
Kikuyu	24.6%	—	—	—
Nandi	25.7%	—	—	—
Swahili	—	53.3%	51.1%	55.5%

435

436

437

438

439

440

441 Table 9 shows that alternation is rare (3–5%), with synthetic and natural values well aligned. Swahili

442 (Natural) alternates most, likely due to shorter, more frequent insertions, whereas Luo tends toward

443 longer insertions and fewer switches.

444

445

446 Table 9: Average alternation rate: percentage of segment boundaries where the language changes.

Source Language	Alternation Rate (%)
Swahili (Synthetic)	4.7%
Luo (Synthetic)	3.6%
Kikuyu (Synthetic)	3.1%
Nandi (Synthetic)	4.1%
Swahili (Natural)	5.3%
Luo (Natural)	2.1%
Average (Synthetic)	3.9%

447

448

449

450

451

452 Across these four dimensions, our model exhibits realistic code-switching structure: it avoids over-

453 insertion, respects language constraints, mirrors natural switch placement, and matches alternation

454 rates—without any hand-crafted rules over language labels. This suggests that structural patterns are

455 implicitly internalized from monolingual segments plus guided diffusion.

456

457

5.2 HUMAN PREFERENCE EVALUATION: FLUENCY AND ACCEPTABILITY

458

459

460

461

462

463

464 To complement automatic metrics, we conduct a large-scale human study assessing the perceived
465 *fluency*, *coherence*, and *realism* of generated code-switched speech—dimensions not fully captured
466 by automated measures. A total of 638 undergraduate participants rated 1,437 utterances sampled
467 across all source languages, with six utterances per listener matched to their linguistic background
468 and at least four independent ratings per utterance.

469

470

471 Participants used a 5-point Likert scale to evaluate:

472

473

474

475

476

- **Fluency:** smoothness and naturalness;
- **Coherence:** semantic consistency and speaker preservation;
- **Realism:** resemblance to naturally occurring multilingual speech.

477

478

479

480

481

482

483

484 Table 10: Average human ratings of code-switched utterances (Likert scale). Natural examples
485 included for reference.

Source Language	Fluency (↑)	Coherence (↑)	Realism (↑)	Std. Dev.
Swahili (Synthetic)	4.1	4.2	4.0	0.42
Luo (Synthetic)	4.1	4.0	4.1	0.46
Kikuyu (Synthetic)	4.2	4.1	4.0	0.44
Nandi (Synthetic)	3.9	4.0	3.9	0.48
Swahili (Natural)	4.6	4.8	4.5	0.42
Luo (Natural)	4.7	4.4	4.6	0.46
Avg. (Synthetic)	4.1	4.05	4.0	0.45

486

487

488

489

490

491

492 Synthetic utterances receive high ratings across all dimensions (≥ 4.0) with low variance across
493 languages. Natural speech scores slightly higher, particularly in realism, but the gap remains modest
494 (typically ≤ 0.5). Overall, listeners perceive the generated speech as fluent, coherent, and plausibly
495 multilingual, supporting the effectiveness of our constraint-guided diffusion approach.

496

497

498

499

500

6 CONCLUSION

501

502

503

504

505

506

507

508

509

510

511 We presented a diffusion-based framework for generating fluent, coherent, and sociolinguistically
512 realistic code-switched speech without relying on parallel code-switched data. By guiding a pre-
513 trained monolingual diffusion prior with differentiable linguistic constraints—including a multilingual514 Table 8: Percentage of foreign segments per ut-
515 terance quarter (Q1: start, Q4: end).

Source	Q1	Q2	Q3	Q4
Swahili (Synthetic)	18.6%	23.5%	26.1%	31.8%
Luo (Synthetic)	22.3%	25.7%	25.1%	26.9%
Kikuyu (Synthetic)	20.8%	22.0%	27.3%	29.9%
Nandi (Synthetic)	19.5%	24.6%	28.4%	27.5%
Swahili (Natural)	16.6%	22.5%	26.1%	34.8%
Luo (Natural)	18.3%	24.7%	23.1%	24.9%

486 language classifier and a contrastive segment encoder—our method performs targeted segment-level
 487 edits while preserving fluency, semantic coherence, and speaker identity.
 488

489 Extensive evaluation across five African languages shows that the system closely matches natural
 490 code-switching behavior in frequency, structure, and temporal placement. It achieves strong segment-
 491 level semantic fidelity (COMET 0.815, LaBSE 0.880) and speaker consistency (EER 6.7%). Human
 492 listeners also rated the generated utterances highly across fluency, coherence, and realism.
 493

494 To our knowledge, this is the first method to enable plug-and-play multilingual infusion within a
 495 single utterance, offering a flexible approach to cross-lingual speech generation in low-resource
 496 settings. Future work will explore richer prosodic control, expansion to additional languages, and
 497 applications to spontaneous conversational speech.
 498

REFERENCES

499 Peter Auer. *Code-switching in conversation: Language, interaction and identity*. Routledge, 1998.
 500

501 Astik Biswas, Emre Yilmaz, Ewald van der Westhuizen, Febe de Wet, and Thomas Niesler. Code-
 502 switched automatic speech recognition in five south african languages. *Computer Speech &*
 503 *Language*, 71:101262, 2022.

504 Yuewen Cao, Songxiang Liu, Xixin Wu, Shiyin Kang, Peng Liu, Zhiyong Wu, Xunying Liu, Dan Su,
 505 Dong Yu, and Helen Meng. Code-switched speech synthesis using bilingual phonetic posteriorgram
 506 with only monolingual corpora. In *ICASSP 2020-2020 IEEE International Conference on Acoustics,*
 507 *Speech and Signal Processing (ICASSP)*, pp. 7619–7623. IEEE, 2020.

508 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
 509 contrastive learning of visual representations. In *International conference on machine learning*, pp.
 510 1597–1607. PMLR, 2020.

512 Jie Chi, Brian Lu, Jason Eisner, Peter Bell, Preethi Jyothi, and Ahmed M Ali. Unsupervised code-
 513 switched text generation from parallel text. In *Proc. INTERSPEECH*, volume 2023, pp. 1419–1423,
 514 2023.

515 Hyungjin Chung, Jonathan Ho, Tim Salimans, Jong Chul Lee, and Diederik P. Kingma. Diffusion
 516 models as plug-and-play priors. In *International Conference on Learning Representations (ICLR)*,
 517 2023. URL <https://openreview.net/forum?id=PlKWVd2yBkY>.

519 Brecht Desplanques, Jenthe Thienpondt, and Kris Demuynck. Ecapa-tdnn: Emphasized chan-
 520 nel attention, propagation and aggregation in tdnn based speaker verification. *arXiv preprint*
 521 *arXiv:2005.07143*, 2020.

523 Paul-Ambroise Duquenne, Hongyu Gong, and Holger Schwenk. Multimodal and multilingual
 524 embeddings for large-scale speech mining. *Advances in Neural Information Processing Systems*,
 525 34:15748–15761, 2021.

526 Fuli Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. Language-agnostic
 527 bert sentence embedding. In *Proceedings of the 60th Annual Meeting of the Association for*
 528 *Computational Linguistics (Volume 1: Long Papers)*, pp. 878–891. Association for Computational
 529 Linguistics, 2022.

531 Bryan Gregorius and Takeshi Okadome. Generating code-switched text from monolingual text with
 532 dependency tree. In *Proceedings of the 20th Annual Workshop of the Australasian Language*
 533 *Technology Association*, pp. 90–97, 2022.

534 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
 535 *neural information processing systems*, 33:6840–6851, 2020.

536 I Hsu, Avik Ray, Shubham Garg, Nanyun Peng, Jing Huang, et al. Code-switched text synthesis in
 537 unseen language pairs. *arXiv preprint arXiv:2305.16724*, 2023.

539 Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. In *IEEE*
 540 *Transactions on Big Data*, volume 7, pp. 535–547. IEEE, 2019.

540 Sehoon Kim, Amir Gholami, Albert Shaw, Nicholas Lee, Karttikeya Mangalam, Jitendra Malik,
 541 Michael W Mahoney, and Kurt Keutzer. Squeezeformer: An efficient transformer for automatic
 542 speech recognition. *Advances in Neural Information Processing Systems*, 35:9361–9373, 2022.
 543

544 Philipp Koehn. Statistical significance tests for machine translation evaluation. In *Proceedings of the
 545 2004 conference on empirical methods in natural language processing*, pp. 388–395, 2004.

546 Carol Myers-Scotton. *Social motivations for code-switching: Evidence from Africa*. Oxford University
 547 Press, 1993.

548

549 Oriol Nieto, Zeyu Jin, Franck Dernoncourt, and Justin Salamon. Efficient spoken language recognition
 550 via multilabel classification. *arXiv preprint arXiv:2306.01945*, 2023.

551

552 Peter Ochieng and Dennis Kaburu. Phonology-guided speech-to-speech translation for african
 553 languages. *Speech Communication*, 174:103287, 2025. ISSN 0167-6393. doi: <https://doi.org/10.1016/j.specom.2025.103287>. URL <https://www.sciencedirect.com/science/article/pii/S0167639325001025>.

554

555 Shana Poplack. Sometimes i'll start a sentence in spanish y termino en espaÑol: toward a typology
 556 of code-switching. *Linguistics*, 18(7-8):581–618, 1980.

557

558 Matt Post. A call for clarity in reporting bleu scores. *arXiv preprint arXiv:1804.08771*, 2018.

559

560 Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever.
 561 Robust speech recognition via large-scale weak supervision. In *International Conference on
 562 Machine Learning*, pp. 28492–28518. PMLR, 2023.

563

564 Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. Comet: A neural framework for mt
 565 evaluation. *arXiv preprint arXiv:2009.09025*, 2020.

566

567 S Sitaram, KR Chandu, SK Rallabandi, and AW Black. A survey of code-switched speech and
 568 language processing. *arXiv preprint arXiv:1904.00784*, 2019.

569

570 Sarah Slabbert and Rosalie Finlayson. A socio-historical overview of codeswitching studies in the
 571 african languages. *South African Journal of African Languages*, 19(1):60–72, 1999.

572

573 Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks.
 574 In *International conference on machine learning*, pp. 6105–6114. PMLR, 2019.

575

576 Ishan Tarunesh, Syamantak Kumar, and Preethi Jyothi. From machine translation to code-switching:
 577 Generating high-quality code-switched text. *arXiv preprint arXiv:2107.06483*, 2021.

578

579 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating
 580 text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.

581

580 7 APPENDIX / SUPPLEMENTAL MATERIAL

582 A DERIVATION OF THE FREE-ENERGY OBJECTIVE

584 We seek to approximate the constrained posterior

$$586 \quad p(x | y) \propto p(x) C(x, y), \quad C(x, y) = c_1(x, y) c_2(x, y),$$

588 where x is an utterance waveform, y is the infusion specification, and c_1, c_2 are soft constraints
 589 (switch plausibility and semantic consistency). Direct inference in $p(x | y)$ is intractable, so we
 590 introduce latent diffusion variables h and consider the joint model $p(x, h) = p(h) p(x | h)$.

591 Starting from the usual variational formulation for $p(x | y) \propto p(x)C(x, y)$, we define the free-energy
 592 functional under a variational distribution $q(x, h)$ as

$$593 \quad F(q) = -\mathbb{E}_{q(x)}[\log p(x)] - \mathbb{E}_{q(x)}[\log C(x, y)] + \mathbb{E}_{q(x)}[\log q(x)]. \quad (9)$$

594 To handle the latent variables, we write
 595

$$596 \log p(x) = \log \int p(x, h) dh = \log \int q(h | x) \frac{p(x, h)}{q(h | x)} dh.$$

598 Applying Jensen's inequality gives the standard variational bound
 599

$$600 \log p(x) \geq \mathbb{E}_{q(h|x)} \left[\log \frac{p(x, h)}{q(h | x)} \right].$$

602 Substituting this bound into equation 9 yields an upper bound on $F(q)$:
 603

$$604 F(q) \leq -\mathbb{E}_{q(x)q(h|x)} \left[\log \frac{p(x, h)}{q(h | x)} \right] + \mathbb{E}_{q(x)} [\log q(x)] - \mathbb{E}_{q(x)} [\log C(x, y)].$$

607 Using $q(x, h) = q(x)q(h | x)$, we can rewrite the first two terms as
 608

$$609 -\mathbb{E}_{q(x)q(h|x)} \left[\log \frac{p(x, h)}{q(h | x)} \right] + \mathbb{E}_{q(x)} [\log q(x)] = \mathbb{E}_{q(x,h)} \left[\log \frac{q(x, h)}{p(x, h)} \right],$$

611 so that the bound takes the familiar free-energy form
 612

$$613 F(q) \leq \mathbb{E}_{q(x,h)} \left[\log \frac{q(x, h)}{p(x, h)} \right] - \mathbb{E}_{q(x)} [\log C(x, y)].$$

615 Motivated by this, we minimize the corresponding free-energy objective stated in Eq. 2 as
 616

$$617 F(q) = \text{KL}(q(x, h) \| p(x, h)) - \mathbb{E}_{q(x)} [\log C(x, y)]. \quad (10)$$

618 The first term encourages samples that are likely under the diffusion prior $p(x, h)$; the second injects
 619 plug-and-play guidance from the soft constraints. In practice we adopt the mode-seeking choice
 620 $q(x) = \delta(x - \eta)$, so that $-\log C(\eta, y)$ appears as an additive guidance penalty inside the reverse
 621 diffusion updates. This connects directly to the sampling algorithm in Algorithm 1.
 622

623 B DDPMs, MODE-SEEKING APPROXIMATION, AND FULL FREE-ENERGY 624 DERIVATION

627 This appendix gives the complete derivation of the constrained free-energy objective used in our
 628 guided diffusion framework. It unifies (i) the DDPM prior, (ii) the variational formulation of the
 629 constrained posterior, (iii) the mode-seeking approximation, and (iv) its reduction to the practical
 630 noise-prediction loss in Eq. 4.
 631

632 B.1 DDPM FORWARD AND REVERSE PROCESSES

633 A denoising diffusion probabilistic model (DDPM) Ho et al. (2020) defines a forward noising process
 634

$$635 q(h = \{x_1, \dots, x_T\} | x_0) = \prod_{t=1}^T q(x_t | x_{t-1}), \quad q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha_t} x_{t-1}, (1 - \alpha_t)I),$$

638 with marginal
 639

$$640 q(x_t | x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t)I), \quad \bar{\alpha}_t = \prod_{s=1}^t \alpha_s.$$

642 The reverse generative process is a Markov chain
 643

$$644 p_\theta(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)), \quad x_T \sim \mathcal{N}(0, I).$$

645 Ho et al. (2020) show that maximizing $\log p_\theta(x_0)$ is equivalent (up to constants) to minimizing the
 646 denoising score-matching loss
 647

$$\mathcal{L}_{\text{DDPM}} = \mathbb{E}_{t, \epsilon_0} [\|\hat{e}_\theta(x_t, t, M) - \epsilon_0\|_2^2], \quad x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_0.$$

648 B.2 CONSTRAINED POSTERIOR AND VARIATIONAL FREE ENERGY
649650 Our target posterior is
651

652
$$p(x | y) \propto p(x) C(x, y), \quad C(x, y) = c_1(x, y) c_2(x, y).$$

653 Since $p(x)$ is defined through the latent diffusion trajectory $h = \{x_1, \dots, x_T\}$, direct inference is
654 intractable, so we use the free-energy objective in Eq. 2 to approximate it.
655656 B.3 MODE-SEEKING APPROXIMATION
657658 Following Chung et al. (2023), we adopt a mode-seeking variational family
659

660
$$q(x) = \delta(x - \eta),$$

661 giving

662
$$q(x, h) = \delta(x - \eta) q(h | \eta).$$

663 Substituting into Eq. 2 yields (up to a constant)
664

665
$$F(\eta, q(h | \eta)) = \text{KL}(q(h | \eta) \| p(h | \eta)) - \log C(\eta, y). \quad (11)$$

666 The KL term is precisely the DDPM variational objective; the second term injects constraint guidance.
667668 Using the forward-diffusion posterior,
669

670
$$q(h | \eta) = \prod_{t=1}^T q(x_t | x_{t-1}, \eta), \quad x_0 = \eta,$$

671 the KL decomposes into per-step terms:
672

673
$$\text{KL}(q(h | \eta) \| p(h | \eta)) = \sum_{t=1}^T \text{KL}(q(x_{t-1} | x_t, \eta) \| p_\theta(x_{t-1} | x_t)).$$

674 Reparameterizing
675

676
$$x_t = \sqrt{\bar{\alpha}_t} \eta + \sqrt{1 - \bar{\alpha}_t} \epsilon_0, \quad \epsilon_0 \sim \mathcal{N}(0, I),$$

677 each term becomes a weighted denoising loss Ho et al. (2020):
678

679
$$\text{KL}(q(x_{t-1} | x_t, \eta) \| p_\theta(x_{t-1} | x_t)) = w_t(\beta) \mathbb{E}_{\epsilon_0} [\|\epsilon_0 - \hat{\epsilon}_\theta(x_t, t, M)\|_2^2].$$

680 Thus
681

682
$$F(\eta) = \sum_{t=1}^T w_t(\beta) \mathbb{E}_{\epsilon_0} [\|\epsilon_0 - \hat{\epsilon}_\theta(x_t, t, M)\|_2^2] - \log C(\eta, y) + \text{const.}$$

683 Using the standard DDPM timestep sampling $t \sim U(1, T)$ and absorbing w_t into the learning rate,
684 we obtain the practical objective used in the main text:
685

686
$$F(\eta) = \mathbb{E}_{t, \epsilon_0} [\|\hat{\epsilon}_\theta(x_t, t, M) - \epsilon_0\|_2^2] - \log C(\eta, y), \quad (12)$$

687 where
688

689
$$x_t = \sqrt{\bar{\alpha}_t} \eta + \sqrt{1 - \bar{\alpha}_t} \epsilon_0.$$

690 This shows that constrained sampling corresponds to standard DDPM reverse updates augmented
691 with the guidance term $-\log C(\eta, y)$, enabling plug-and-play code-switching without retraining the
692 diffusion prior.
693

702 **C SOFT TOP- M RETRIEVAL WITH TEMPERATURE ANNEALING**
 703
 704
 705

706 This appendix details an optional variant of the semantic retrieval procedure in Section 3.2. Instead
 707 of selecting a single top-1 nearest neighbor at each denoising step, we form a soft mixture over the
 708 top- M candidates. This stabilizes early reverse-diffusion updates when x_t is still highly noisy and
 709 the encoder embeddings may be unreliable.

710 **Candidate Retrieval and Similarity Logits.** Given a source segment

$$711 \quad q = f_{\text{enc}}(s_{x_j}^{(i^*)}),$$

713 we query the FAISS index \mathcal{D} (restricted to the infusion-eligible language set $\mathcal{S}_{\text{inf}}(y)$) and obtain the
 714 top- M candidates:

$$715 \quad \mathcal{C}_M = \{s_{y_k}^{(m_1)}, \dots, s_{y_k}^{(m_M)}\}.$$

716 For each candidate we compute cosine-similarity logits

$$718 \quad z_r = \text{Sim}(q, f_{\text{enc}}(s_{y_k}^{(m_r)})), \quad r = 1, \dots, M.$$

720 **Soft Retrieval Distribution.** A tempered softmax converts the logits into a probability distribution:

$$722 \quad \pi_r(t) = \frac{\exp(z_r/\tau_t)}{\sum_{u=1}^M \exp(z_u/\tau_t)}.$$

724 Large temperatures τ_t yield diffuse mixtures (exploration), while $\tau_t \rightarrow 0$ collapses the distribution to
 725 the best candidate.

727 **Mixture-Based Segment Construction.** The retrieved segment is the convex combination

$$729 \quad s^*(t) = \sum_{r=1}^M \pi_r(t) s_{y_k}^{(m_r)}.$$

732 As $\tau_t \rightarrow 0$, the mixture degenerates to the hard top-1 candidate:

$$733 \quad s^*(t) \longrightarrow s_{y_k}^{(m^*)}.$$

736 **Temperature Annealing Schedule.** We anneal τ_t across reverse-diffusion steps so the model explores
 737 early and commits later. A simple schedule is

$$739 \quad \tau_t = \tau_{\min} + (\tau_{\max} - \tau_{\min}) \left(1 - \frac{t}{T}\right)^\kappa, \quad \kappa \in [2, 4], \quad \tau_{\min} \leq \tau_t \leq \tau_{\max},$$

741 where $\tau_{\max} \approx 1.0\text{--}2.0$ and τ_{\min} (e.g. 10^{-3}) prevents numerical instability. Here, $\tau_t \approx \tau_{\max}$ when
 742 $t \approx 1$ (early noisy steps, more exploration) and decays toward τ_{\min} as $t \rightarrow T$ (later steps, sharper
 743 selection).

745 **Gradient Flow.** Gradients propagate through the mixture weights $\pi_r(t)$, through the dependence of
 746 the embeddings $f_{\text{enc}}(s_{y_k}^{(m_r)})$ and $f_{\text{enc}}(s_{x_j}^{(i^*)})$ on the clean sample η , and through the blend-and-write-
 747 back update in Section 3.2. Both the FAISS index and the encoder *parameters* remain frozen; only
 748 the clean sample η receives updates from \mathcal{L}_{c_2} .

749 **When Soft Top- M Helps.** The soft retrieval variant is especially useful when:

751

- 752 • languages in $\mathcal{S}_{\text{inf}}(y)$ are phonetically similar (ambiguous nearest neighbours),
- 753 • segments are short (50–100 ms), making embeddings noise-sensitive,
- 754 • early DDPM steps ($t \approx T$) are dominated by noise.

755 As $\tau_t \downarrow 0$, the method reduces to the deterministic top-1 retrieval in the main text.

756 **D PROSODIC NORMALIZATION AND TOLERANCE PARAMETERS**
 757

758 For the prosodic compatibility terms in Section 3.2, we normalize candidate durations using simple
 759 speech-rate and prosody proxies. Let R_{x_j}, P_{x_j} and R_{y_k}, P_{y_k} denote speech-rate and prosodic
 760 statistics (e.g., syllables/s, median F0 or energy) for the host utterance x_j and the candidate utterance
 761 y_k , respectively. We define a scale factor

$$762 \quad S_{\text{ratio}} = \sqrt{\frac{R_{x_j}}{R_{y_k}} \cdot \frac{P_{x_j}}{P_{y_k}}}$$

$$763$$

$$764$$

$$765$$

766 and an expected candidate duration

$$767 \quad \hat{d} = d_{x_j}^{(i^*)} S_{\text{ratio}},$$

$$768$$

769 where $d_{x_j}^{(i^*)}$ is the duration of the source segment chosen for infusion.

770 The duration tolerance parameter $\lambda_d \in [0, 1)$ controls the allowable relative deviation from \hat{d} . In the
 771 main text, we use the hinge penalty

$$772$$

$$773 \quad \mathcal{L}_{\text{dur}} = \max\left(0, \frac{|d_{y_k}^{(m^*)} - \hat{d}| - \lambda_d \hat{d}}{\hat{d}}\right),$$

$$774$$

$$775$$

776 which becomes zero when the candidate duration lies within a $(1 \pm \lambda_d)$ band around \hat{d} .

777 For onset alignment we define

$$778 \quad \Delta\tau = |\bar{d}_{x_j} - \bar{d}_{y_k}|,$$

$$779$$

780 where \bar{d}_{x_j} and \bar{d}_{y_k} are the mean segment durations in x_j and y_k , respectively. The onset penalty

$$781 \quad \mathcal{L}_{\text{on}} = \max\left(0, \frac{|\mathcal{O}_{y_k}^{(m^*)} - \mathcal{O}_{x_j}^{(i^*)}| - \Delta\tau}{\Delta\tau}\right)$$

$$782$$

$$783$$

784 is therefore zero when the candidate onset falls within a tolerance window of width $\Delta\tau$ around the
 785 host onset. In our experiments we set λ_d and any additional scaling of $\Delta\tau$ via development tuning on
 786 a held-out validation set

$$787$$

788 **E NEIGHBORHOOD DEFINITION FOR CONTEXTUAL CONSISTENCY**
 789

790 For the contextual loss in Section 3.2, we define the neighborhood $\mathcal{N}(i^*)$ of the edited segment $s_{x_j}^{(i^*)}$
 791 as a fixed-radius window over adjacent host segments. Let x_j be segmented into n spans $\{s_{x_j}^{(i)}\}_{i=1}^n$,
 792 and let i^* be the index selected for infusion. For a window radius $R \in \mathbb{N}$ (typically $R = 1$ or $R = 2$),
 793 the neighborhood is

$$794$$

$$795 \quad \mathcal{N}(i^*) = \{s_{x_j}^{(i)} : \max(1, i^* - R) \leq i \leq \min(n, i^* + R), i \neq i^*\}.$$

$$796$$

$$797$$

798 This definition automatically excludes segments outside the valid range $[1, n]$ and captures local
 799 prosodic and semantic context around the infusion site.

$$800$$

801 **Special case ($R = 1$).** For immediate left/right neighbors, the above reduces to the standard
 802 adjacent-neighbor definition:

$$803 \quad \mathcal{N}(i^*) = \begin{cases} \{s_{x_j}^{(i^*+1)}\}, & i^* = 1, \\ \{s_{x_j}^{(i^*-1)}, s_{x_j}^{(i^*+1)}\}, & 1 < i^* < n, \\ \{s_{x_j}^{(i^*-1)}\}, & i^* = n. \end{cases}$$

$$804$$

$$805$$

$$806$$

$$807$$

808 The general-radius formulation allows broader contextual windows when desired, while the $R = 1$
 809 instance recovers the conventional adjacent-segment neighborhood.

810 F FULL INFERENCE ALGORITHM AND GRADIENT UPDATE SCHEDULE 811

812 This appendix expands on Section 3.3 of the main text, providing full details on the inference
813 algorithm, segment-wise gradient updates, constraint scheduling, and hyperparameter tuning for
814 code-switched speech generation. We elaborate on the free-energy formulation in Eq. 4 and describe
815 the iterative refinement strategy that supports semantically and prosodically aligned code-switching.
816

817 F.1 INFERENCE OBJECTIVE 818

819 Recall the guided per-step loss from Eq. 7:
820

$$821 \mathcal{L}_{\text{step}}(t) = \|\hat{\epsilon}_\theta(x_t, t, M) - \epsilon_0\|_2^2 + \lambda_1 \mathcal{L}_{c_1} + \lambda_2(t) \mathcal{L}_{c_2},$$

822 where $x_t = \sqrt{\bar{\alpha}_t} \eta + \sqrt{1 - \bar{\alpha}_t} \epsilon_0$ and $\epsilon_0 \sim \mathcal{N}(0, I)$. The DDPM parameters θ are pretrained and
823 frozen; we optimize only the clean sample η . Averaging over timesteps and noise draws yields the
824 overall guided objective
825

$$826 F(\eta) = \mathbb{E}_{t, \epsilon_0} [\|\hat{\epsilon}_\theta(x_t, t, M) - \epsilon_0\|_2^2 + \lambda_1 \mathcal{L}_{c_1} + \lambda_2(t) \mathcal{L}_{c_2}], \quad (13)$$

827 which instantiates the free-energy form in Eq. 4 with explicit weights on the constraint terms. Here,
828 \mathcal{L}_{c_1} and \mathcal{L}_{c_2} are defined in Sections 3.1 and 3.2, respectively.
829

830 F.2 DYNAMIC SCHEDULING FOR INFUSION 831

832 To control when and how strongly foreign segments are introduced, we define two schedules that
833 match the discussion in Sec. 3.3.
834

835 **Time-dependent blending coefficient.** We use a ramp $\rho_t \in [0, 1]$ to blend the retrieved foreign
836 segment into the host segment (see Sec. 3.2):
837

$$838 s_{x_j}^{(i^*)} \leftarrow (1 - \rho_t) s_{x_j}^{(i^*)} + \rho_t s^*.$$

840 A simple schedule is
841

$$842 \rho_t = 1 - \exp\left(-\frac{T - t}{\beta T}\right), \quad (14)$$

844 with $\beta \in (0, 1)$ (we use $\beta = 0.25$). This ensures that early reverse-diffusion steps (t near T) preserve
845 monolingual structure, while foreign infusion gradually intensifies as t decreases and the acoustic
846 structure stabilizes.
847

848 **Constraint weight ramp-up.** We factor the guidance weight as $\lambda_2(t) = \lambda_2 w(t)$, where
849

$$850 w(t) = \frac{t}{T}, \quad (15)$$

852 so that \mathcal{L}_{c_2} is suppressed when x_t is still highly corrupted and only becomes influential once a coarse
853 waveform has formed. In practice, this reduces the risk of semantically misaligned edits at very noisy
854 timesteps.
855

856 F.3 CONSTRAINT WEIGHT SELECTION 857

858 We select λ_1 and the base λ_2 via Gaussian process-based Bayesian optimization over the range
859 $[0.1, 5.0]$, using a held-out validation set. The objective combines multiple evaluation metrics:
860

- **Semantic fidelity:** COMET, BERTScore;
- **Prosodic alignment:** onset and duration deviation at switch boundaries;
- **Speaker consistency:** cosine similarity using ECAPA-TDNN embeddings.

863 The optimal weights used in our experiments are $\lambda_1 = 0.35$ and $\lambda_2 = 0.65$.

864
865
866
867
868
869

870 At each timestep t_i , we update the clean sample η using gradients that are localized to a single
 871 segment of the utterance. Let $\tilde{x}_{t_i} = \text{Segment}(x_{t_i}, L)$ denote a subwindow of x_{t_i} of length L centred
 872 on the segment index i^* selected by c_1 . The guided loss at step t_i is

873
874
875
876
877
878
879
880
881
882

$$883 \quad \mathcal{L}_{\text{step}}(t_i) = \|\epsilon_0 - \hat{\epsilon}_\theta(\tilde{x}_{t_i}, t_i, M)\|_2^2 + \lambda_1 \mathcal{L}_{c_1} + \lambda_2(t_i) \mathcal{L}_{c_2},$$

884
885
886
887
888
889
890
891
892

893 and we take a gradient step on η :

894
895
896
897
898
899
900
901
902
903

$$904 \quad \eta \leftarrow \eta - \lambda_\eta \nabla_\eta \mathcal{L}_{\text{step}}(t_i), \quad (16)$$

905
906
907
908
909
910
911
912
913
914

915 with a small learning rate λ_η (we use $\lambda_\eta = 0.05$). Because \tilde{x}_{t_i} depends on η only through the
 916 selected span, this update predominantly affects a single localized region of the utterance at each step.
 917 Over the course of the reverse trajectory, different segments are selected, enabling diverse foreign
 substitutions while preserving global fluency and speaker identity.

918 F.5 FULL INFERENCE ALGORITHM
919920
921 **Algorithm 1** Guided DDPM Inference for Code-Switched Speech

922 **Require:** Frozen DDPM denoiser $\hat{\epsilon}_\theta$, frozen LID classifier f_{cl} , frozen multilingual encoder
923 f_{enc} , FAISS-based vector database \mathcal{D} , infusion specification y (incl. \mathcal{S}_{inf} , π_{switch}), host Mel
924 spectrogram M , noise schedule $\{\bar{\alpha}_t\}_{t=1}^T$, step sizes $\{\gamma_t\}_{t=1}^T$, guidance weights $\lambda_1, \lambda_2(t)$, blend
925 ramp ρ_t

926 1: Initialize $\eta \sim \mathcal{N}(0, I)$ {initial guess for the clean sample}

927 2: **for** $t = T, T-1, \dots, 1$ **do**

928 3: // 1. Segment and compute LID-based controller c_1

929 4: Segment η into spans $\{s^{(i)}\}_{i=1}^n$

930 5: For each span $s^{(i)}$: compute $p(\ell | s^{(i)}) = f_{\text{cl}}(s^{(i)})_\ell$ over $\ell \in \{\ell_{\text{mono}}\} \cup \mathcal{S}_{\text{inf}}(y)$

931 6: Compute foreignness scores $u^{(i)} = \sum_{\ell \in \mathcal{S}_{\text{inf}}(y)} p(\ell | s^{(i)})$ and $P_{\text{foreign}}(\eta) = \frac{1}{n} \sum_i u^{(i)}$ (Eq. 5)

932 7: Compute global controller $c_1(\eta, y)$ and local gates $g^{(i)}$ as in Sec. 3.1

933 8: Set $\hat{\mathcal{L}}_{c_1}(\eta, y) = -\log c_1(\eta, y)$ (Eq. 6)

934 9: Choose single span to edit: $i^* = \arg \max_i g^{(i)}$

935 10: // 2. Retrieve and score foreign segment (c_2)

936 11: Encode query $q = f_{\text{enc}}(s^{(i^*)})$

937 12: Query \mathcal{D} restricted to $\ell(s_{y_k}^{(m)}) \in \mathcal{S}_{\text{inf}}(y)$ and select

938
$$m^* = \arg \max_{\substack{s_{y_k}^{(m)} \in \mathcal{D} \\ \ell(s_{y_k}^{(m)}) \in \mathcal{S}_{\text{inf}}(y)}} \text{Sim}(q, f_{\text{enc}}(s_{y_k}^{(m)})), \quad s^* = s_{y_k}^{(m^*)}.$$

939 13: Compute prosody-aware penalties $\mathcal{L}_{\text{dur}}, \mathcal{L}_{\text{on}}$ (Sec. 3.2, App. D)

940 14: Compute semantic and contextual losses $\mathcal{L}_{\text{sem}}, \mathcal{L}_{\text{ctx}}$ (Sec. 3.2)

941 15: Form infusion loss

942
$$\mathcal{L}_{c_2} = g^{(i^*)} [\alpha_{\text{sem}} \mathcal{L}_{\text{sem}} + \alpha_{\text{ctx}} \mathcal{L}_{\text{ctx}} + \alpha_{\text{pro}} (\mathcal{L}_{\text{dur}} + \mathcal{L}_{\text{on}})]$$

943 16: // 3. Blend-and-write-back in the clean domain

944 17: $s^{(i^*)} \leftarrow (1 - \rho_t) s^{(i^*)} + \rho_t s^*$ (update span in η)

945 18: Reassemble η from updated spans $\{s^{(i)}\}$

946 19: // 4. DDPM denoising + guided gradient step

947 20: Sample $\epsilon_0 \sim \mathcal{N}(0, I)$ and set $x_t = \sqrt{\bar{\alpha}_t} \eta + \sqrt{1 - \bar{\alpha}_t} \epsilon_0$

948 21: Predict noise: $\hat{\epsilon}_\theta \leftarrow \hat{\epsilon}_\theta(x_t, t, M)$

949 22: Define step loss

950
$$\mathcal{L}_{\text{step}}(t) = \|\hat{\epsilon}_\theta - \epsilon_0\|_2^2 + \lambda_1 \mathcal{L}_{c_1} + \lambda_2(t) \mathcal{L}_{c_2}$$

951 23: Update clean sample:

952
$$\eta \leftarrow \eta - \gamma_t \nabla_\eta \mathcal{L}_{\text{step}}(t)$$

953 24: **end for**

954 25: **return** η as the code-switched waveform x_0

955

956

957

G IDENTITY REFINEMENT AND SPEAKER HARMONIZATION DETAILS

958

959

960

961 Although the code-switched utterance x is semantically coherent after guided generation, we observe
962 that local segment-level gradients and cross-lingual substitutions can introduce inconsistencies in
963 voice quality, prosody, or timbre. To address this, we add a short identity-harmonization pass using
964 the *same pretrained, frozen* diffusion model as in the main sampler.

965

966

967

968 Given a monolingual reference utterance x_{mono} and the current code-switched sample x , we first
969 extract a global speaker/prosody descriptor

970

971

$$\phi_{\text{spk}} = \text{ECAPA}(x_{\text{mono}}), \quad \phi_{\text{mel}} = \text{MelStats}(x),$$

972 where ECAPA(\cdot) is a frozen ECAPA-TDNN encoder and MelStats(\cdot) denotes global statistics
 973 (mean and variance) of log-Mel features. We bundle these into a single conditioning vector
 974

$$\phi_{\text{target}} = g(\phi_{\text{spk}}, \phi_{\text{mel}}),$$

975 implemented as a small linear projection, and run a shallow denoising refinement
 976

$$x_{\text{final}} = \text{DDPM}_{\text{refine}}(x \mid \phi_{\text{target}}), \quad (17)$$

977 for $T_{\text{ref}} = 150$ late diffusion steps with a low-noise schedule. During this refinement we disable
 978 segment edits by setting $\lambda_2(t) = 0$ and keep the LID-based rate term weak (small λ_1), so semantics
 979 and the code-switch pattern are preserved.

980 Empirically, this post-hoc refinement corrects subtle inconsistencies without altering content. In
 981 particular, it improves:

- 982 • **Timbre smoothing** — reduces artifacts from mismatched vocal-tract characteristics across
 983 segments;
- 984 • **Prosodic coherence** — better alignment of pitch and rhythm across switch boundaries;
- 985 • **Voice uniformity** — the utterance sounds more consistently like a single speaker.

986 We also experimented with adding an explicit speaker-consistency loss $\mathcal{L}_{\text{id}} = 1 -$
 987 $\cos(\text{ECAPA}(x), \text{ECAPA}(x_{\text{mono}}))$ to the guided objective in Eq. 7. However, this often led to
 988 unstable behavior and degraded convergence due to conflicts with semantic and timing objectives.
 989 In contrast, the post-hoc harmonization pass offered better control, computational simplicity, and
 990 training stability, while achieving comparable or better speaker-consistency scores.

991 H TOOLS AND RESOURCES

992 **Multilingual Language Classifier f_{cl} .** We adopt the LECAPAT architecture Nieto et al. (2023), a
 993 lightweight variant of ECAPA-TDNN Desplanques et al. (2020), as our multilingual segment-level
 994 language classifier f_{cl} . The classifier takes log-Mel spectrograms (64 bins, 25 ms window, 10 ms hop,
 995 64 ms FFT) extracted from 24 kHz audio and predicts language identity for each segment.

996 The model is trained with cross-entropy over five languages for 50 epochs using Adam ($\text{lr} = 10^{-4}$,
 997 batch size 64), a 10% validation split, and early stopping (patience: 5). No data augmentation was
 998 used. On a single NVIDIA A100 GPU, the classifier achieves 92.4% average validation accuracy.
 999 During DCSM inference, f_{cl} is frozen and used only to compute the foreignness scores defined in
 1000 §3.1.

1001 **Multilingual Segment Encoder f_s .** The multilingual encoder f_s maps speech segments across
 1002 languages into a shared latent space using a SimCLR-style contrastive objective Chen et al. (2020).
 1003 Positive pairs (same-language segments) are drawn closer in embedding space, while negative pairs
 1004 (cross-language) are pushed apart. To improve robustness, 50% of training segments are augmented
 1005 with Gaussian noise. The encoder architecture includes a 1D convolutional frontend (256 filters,
 1006 kernel size 16, stride 8), followed by an EfficientNet-B0 Tan & Le (2019) backbone and global
 1007 max pooling. A projection head maps representations to a 720-dim contrastive space. The model
 1008 was trained for 1M steps using AdamW ($\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$, batch size 512). Only
 1009 EfficientNet embeddings are used at inference.

1010 **Pre-trained Diffusion Model (SegUniDiff).** For speech generation, we use the Segment-Aware
 1011 Unified Diffusion Model (SegUniDiff) Ochieng & Kaburu (2025), which synthesizes code-switched
 1012 utterances from paired segments (s_{x_i}, s_{y_k}) via a denoising diffusion process. Each model is trained
 1013 per language pair, conditioned on Mel-spectrograms to capture acoustic context. We refer to Ochieng
 1014 & Kaburu (2025) for architectural and training specifics.

1015 **Machine Translation and ASR Models.** To support automatic evaluation, we constructed parallel
 1016 corpora by manually aligning semantically equivalent sentences across all language pairs in our
 1017 dataset. These were used to train Transformer-base machine translation (MT) models. For automatic
 1018 speech recognition (ASR), we trained five language-specific models: Squeezeformer Kim et al. (2022)
 1019 for Nandi, Luo, and Kikuyu, and Whisper-small Radford et al. (2023) for Swahili and English.

1026

1027 Table 11: Parallel MT datasets with SacreBLEU scores and ASR performance by language.

Language Pair	Paired Sentences	SacreBLEU (\uparrow)	Language (ASR)	WER (%)
Luo–Nandi	1.76M	32.2	Luo	14.2
Luo–Kikuyu	1.18M	31.8	Nandi	13.6
Nandi–Kikuyu	1.32M	27.3	Kikuyu	14.4
Kikuyu–Swahili	1.29M	30.4	Swahili	9.8
Kikuyu–English	1.71M	24.9	English	5.3
Swahili–English	1.52M	25.4	—	—
Luo–Swahili	1.43M	27.4	—	—
Luo–English	1.34M	28.1	—	—
Nandi–Swahili	1.44M	27.1	—	—
Nandi–English	1.37M	28.6	—	—

1035

1036

1037 I EVALUATION METRIC DETAILS AND RESAMPLING PROTOCOL

1038 To evaluate semantic and linguistic fidelity of generated code-switched utterances, we use:

- **SacreBLEU** Post (2018): Measures n-gram overlap with detokenization invariance.
- **BERTScore** Zhang et al. (2019): Captures contextual similarity using pre-trained transformer embeddings.
- **COMET** Rei et al. (2020): A learned metric trained on human judgments of adequacy and fluency.
- **LaBSE Similarity** Feng et al. (2022): Cosine similarity between sentence embeddings from multilingual BERT, used on English translations to assess discourse-level alignment.

1050

1051 J RESAMPLING METHOD FOR SEGMENT-LEVEL EVALUATION

1052

1053 To compute 95% confidence intervals for each metric, we adopt a bootstrap resampling procedure Koehn (2004). For each source language:

1. We construct 100 test sets of 700 segment pairs $(s_x^{(k)}, s_{x_c}^{(k)})$ sampled from the full evaluation pool.
2. For each test set, we perform 1,000 bootstrap iterations by sampling with replacement.
3. In each iteration, we concatenate all reference translations (from $s_x^{(k)}$) into one string and all hypothesis translations (from $s_{x_c}^{(k)}$) into another.
4. We compute SacreBLEU, BERTScore, COMET, and LaBSE similarity between these concatenated sequences.
5. We report the 95% confidence interval as the range between the 2.5th and 97.5th percentiles of the resulting score distributions.

1065 This procedure ensures statistically stable estimates across a diverse evaluation population, capturing 1066 both ASR/MT variability and segment-level diversity.

1067

1068 K RESAMPLING PROCEDURE FOR UTTERANCE-LEVEL EVALUATION

1069

1070 To compute confidence intervals for utterance-level metrics, we follow this bootstrap-based procedure:

1. Construct 100 test sets per language, each with 700 utterance pairs (x, x_r) or (x, x_m) .
2. Transcribe all utterances using language-specific ASR systems.
3. Translate into English (if not already monolingual).
4. Perform 1,000 bootstrap iterations:
 - Sample 700 utterance pairs with replacement.
 - Concatenate references and hypotheses into long sequences.
 - Compute SacreBLEU, BERTScore, COMET, and LaBSE similarity.
5. Report mean and 95% CI from the score distribution (2.5th–97.5th percentiles).

1080 K.1 TOLERANCE SELECTION FOR CROSS-LINGUAL SEGMENT SUBSTITUTION
10811082 To ensure rhythmic and temporal alignment during code-switched segment substitution, we adopt a
1083 data-driven strategy for selecting the tolerance parameter λ . This parameter governs the allowable
1084 deviation between the duration of a monolingual segment and that of a candidate segment drawn
1085 from an infusion language.1086 For each pair of languages involved in substitution—where one provides the *monolingual segment*
1087 and the other contributes the *infused segment*—we compute a base tolerance as the normalized
1088 difference in average segment durations:
1089

1090
$$\lambda_{\text{base}}(\ell_x, \ell_y) = \frac{|\bar{d}_{\ell_x} - \bar{d}_{\ell_y}|}{\bar{d}_{\ell_x}}, \quad (18)$$

1091

1092 where \bar{d}_{ℓ_x} and \bar{d}_{ℓ_y} denote the average segment durations (in seconds) for the monolingual and infusion
1093 languages, respectively. This base ratio captures prosodic variation and relative speaking rates
1094 between languages.
10951096 The final tolerance is then defined as:
1097

1098
$$\lambda(\ell_x, \ell_y) = \max(\lambda_{\text{base}}(\ell_x, \ell_y) + \epsilon, \lambda_{\min}), \quad (19)$$

1099

1100 where ϵ is a fixed safety margin (set to 0.05), and λ_{\min} is a lower bound (set to 0.1) to prevent
1101 over-constraining substitutions in closely matched language pairs. This formulation allows λ to scale
1102 naturally with inter-language temporal divergence, while preserving a minimal tolerance window
1103 across all combinations.
11041105 Table 12: Computed λ values for Swahili (ℓ_x) as the monolingual language. Average durations are in
1106 seconds.
1107

Infusion Language ℓ_y	Avg. Duration d_{ℓ_y} (s)	λ_{base}	Final λ
Luo	16.5	0.0855	0.1355
Nandi	17.1	0.1250	0.1750
Kikuyu	15.6	0.0263	0.1000
English	15.0	0.0132	0.1000
Swahili average segment duration: $\bar{d}_{\ell_x} = 15.2$ seconds			

1113 In practice, these empirically derived λ values led to high substitution success rates and prosodically
1114 natural code-switched utterances across language pairs. The approach enabled rhythm-preserving
1115 segment replacement while maintaining tight control over misaligned insertions.
11161117 L ABLATION
11181119 L.1 EFFECT OF REMOVING THE LANGUAGE CLASSIFIER CONSTRAINT
11201121 In this experiment, we evaluate the impact of removing the first constraint $c_1(x, y)$, which guides
1122 language identification and determines the location of foreign segment insertions. During inference,
1123 we modify the loss function by omitting the classifier-related term, resulting in the following objective:
1124

1125
$$F = \arg \min_{\theta} \mathbb{E}_{t \sim U(2, T)} \left[\left\| \hat{\epsilon}_{\theta}(x_t, t, M) - \epsilon_0 \right\|_2^2 \right] + \mathcal{L}_{c_2}. \quad (20)$$

1126

1127 We analyze the behavior of this classifier-free variant along three key dimensions of code-switching:
1128 **(i)** the frequency of switching, **(ii)** alternation points, and **(iii)** subjective fluency, coherence, and
1129 realism as rated by human evaluators. From the full set of 8,500 generated code-switched utterances,
1130 we randomly sample 2,000 utterances per source language and follow the evaluation procedures
1131 described in Sections 5.1 and 5.2.
11321133 Table 13 presents a comparison between the full model and the variant without the language classifier
constraint. Removing \mathcal{L}_{c_1} results in a substantial increase in code-switching frequency—from 4.33%

1134 to 18.3%—and a corresponding spike in alternation rate—from 3.88% to 17.9%. These shifts indicate
 1135 that, without a mechanism to regulate switch locations, the model overproduces foreign segments
 1136 and places them erratically throughout the utterance.

1137 This overgeneration directly impacts speech naturalness. Human ratings reveal a marked decline
 1138 in fluency (from 4.1 to 2.7), coherence (from 4.05 to 3.3), and realism (from 4.0 to 2.6). These
 1139 results underscore the importance of the classifier constraint in producing linguistically appropriate,
 1140 contextually coherent, and perceptually natural code-switched speech.

1141
 1142
 1143 Table 13: Comparison of code-switching behavior between the full model and the variant without the
 1144 language classifier constraint \mathcal{L}_{c_1} .

Metric	Full Model	Without \mathcal{L}_{c_1}	Difference (Δ)
Avg. Code-Switching Frequency	4.33%	18.3%	+13.97%
Avg. Alternation Rate	3.88%	17.9%	+14.02%
Avg. Fluency (Human)	4.1	2.7	-1.4
Avg. Coherence (Human)	4.05	3.3	-0.75
Avg. Realism (Human)	4.0	2.6	-1.4

1150 1151 L.2 EFFECT OF REMOVING TEMPORAL ALIGNMENT AND ONSET CONSTRAINTS 1152

1153 In this experiment, we assess the impact of removing the duration and onset components of the
 1154 constraint loss \mathcal{L}_{c_2} , which enforce prosodic alignment between the inserted foreign segment and the
 1155 original monolingual utterance. These constraints ensure that inserted segments match the expected
 1156 duration and start at a position consistent with the rhythm and flow of the host utterance, thereby
 1157 preserving fluency and naturalness.

1158 To isolate their contribution, we exclude both terms by setting $\gamma = 0$, resulting in a simplified
 1159 constraint loss:

$$\mathcal{L}_{c_2} = \alpha \cdot \mathcal{L}_{\text{semantic}} + \beta \cdot \mathcal{L}_{\text{context}}.$$

1160 This ablation allows the model to insert segments of arbitrary duration and onset without explicit
 1161 prosodic guidance. We generate a total of 8,500 code-switched utterances and evaluate both segment-
 1162 level and utterance-level quality using the procedures described in Sections ?? and ??.

1163 Table ?? summarizes the average performance across all synthetic languages, with and without the
 1164 duration/onset constraints.

1165
 1166 Table 14: Impact of removing duration and onset constraints on segment- and utterance-level
 1167 evaluation metrics. Scores are reported as mean values with 95% confidence intervals (CI) based on
 1168 1,000 bootstrap samples.

Level	Metric	Avg. With Duration/Onset	Avg. Without Duration/Onset
Segment	SacreBLEU (\uparrow)	36.4 [35.2, 37.4]	34.6 [33.3, 35.7]
	BERTScore (\uparrow)	0.807 [0.804, 0.810]	0.796 [0.792, 0.800]
	COMET (\uparrow)	0.815 [0.807, 0.823]	0.790 [0.781, 0.799]
	LaBSE Similarity (\uparrow)	0.880 [0.873, 0.886]	0.868 [0.860, 0.874]
Utterance	SacreBLEU (\uparrow)	34.5 [33.5, 35.6]	32.8 [31.5, 33.9]
	BERTScore (\uparrow)	0.755 [0.751, 0.759]	0.743 [0.738, 0.748]
	COMET (\uparrow)	0.653 [0.645, 0.661]	0.624 [0.614, 0.634]
	LaBSE Similarity (\uparrow)	0.874 [0.870, 0.878]	0.860 [0.854, 0.867]

1178 Table 15: Human evaluation scores comparing the full model with the variant without duration/onset
 1179 constraints. Ratings are on a 5-point Likert scale.

Metric	Full Model	Without Duration/Onset	Difference (Δ)
Avg. Fluency (Human)	4.1	3.1	-1.0
Avg. Coherence (Human)	4.05	3.3	-0.75
Avg. Realism (Human)	4.0	2.4	-1.6

1184 Quantitative results in Table 14 show consistent declines in both segment- and utterance-level metrics
 1185 across all evaluation measures. While the degradation in segment-level scores is modest (e.g., -1.8
 1186 SacreBLEU, -0.019 COMET), utterance-level metrics are more sensitive to prosodic disruptions,
 1187 with COMET and BERTScore dropping by 0.029 and 0.012, respectively. These drops suggest that

1188 even minor misalignments in duration or onset can propagate across an utterance, leading to broader
 1189 semantic and rhythmic incoherence.

1190
 1191 Human evaluations (Table 15) further confirm these effects. Fluency and realism drop significantly
 1192 (by -1.0 and -1.6 points, respectively), with listeners noting more jarring transitions and unnatural
 1193 pacing. Although semantic coherence is partially preserved (-0.75), the lack of prosodic control leads
 1194 to degraded overall acceptability.

1195 Together, these results highlight the critical role of timing constraints in producing fluent and natural-
 1196 sounding code-switched speech. Their removal leads to audible temporal mismatches, underscoring
 1197 the need to model prosodic structure alongside semantics and context

1200 M RELATED WORK

1201
 1202 **Code-switching** is a well-documented linguistic phenomenon in multilingual communities, particu-
 1203 larly across Africa, where speakers frequently alternate between local vernaculars and national or
 1204 international languages such as English or Swahili. Foundational work by Slabbert & Finlayson
 1205 (1999) and Myers-Scotton (1993) highlighted code-switching as a communicative strategy influenced
 1206 by identity, context, and pragmatics. Poplack (1980) and Auer (1998) further explored structural
 1207 patterns and conversational dynamics, establishing typologies of alternation, insertion, and congruent
 1208 lexicalization. These studies underscore the naturalness and linguistic richness of code-switching in
 1209 African speech.

1210 Despite its sociolinguistic prominence, *code-switching has been underrepresented in computational*
 1211 *speech research*, largely due to the lack of annotated corpora and standardized tools. While progress
 1212 has been made in *code-switched text generation* using statistical or neural methods (Tarunesh et al.,
 1213 2021; Gregorius & Okadome, 2022; Chi et al., 2023), the *speech modality* remains significantly
 1214 underexplored.

1215 The most notable contribution to *code-switched speech synthesis* is by Cao et al. (2020), who proposed
 1216 a bilingual phonetic posteriorgram-based model that combines monolingual speech corpora to gener-
 1217 ate mixed-language speech. However, their method lacks explicit semantic or contextual alignment
 1218 and does not account for speaker consistency or natural prosodic transitions across languages.

1219 In contrast, our work introduces a *diffusion-based framework* that synthesizes code-switched speech
 1220 by minimally editing monolingual utterances. We incorporate *linguistic constraints*—a pre-trained lan-
 1221 guage classifier for soft switch control and a multilingual encoder for semantic segment matching—to
 1222 guide the generation process. Additionally, we address *speaker identity harmonization* by introducing
 1223 a refinement step based on acoustic conditioning.

1224 *To the best of our knowledge, this is the first work that enables the infusion of multiple foreign*
 1225 *languages within a single utterance, allowing for rich, naturalistic multilingual code-switching*
 1226 *patterns.* This represents a significant advancement toward realistic speech generation in low-resource
 1227 multilingual settings.

1228 1229 N LIMITATIONS

1230 Our proposed framework for controlled code-switched speech generation has demonstrated strong
 1231 quantitative and human evaluation performance. However, several limitations remain:

1232 Mismatch Between Synthesized and Natural Speech The generated utterances, while fluent and
 1233 semantically faithful, are synthesized from noise and do not inherit the rich socio-pragmatic cues,
 1234 emotional tone, or discourse-driven switching patterns present in natural conversations. This limits
 1235 the realism of certain paralinguistic features such as emphasis, hesitation, or spontaneous repairs.

1236 No Parallel Code-Switched Supervision The model is trained entirely on monolingual utterances
 1237 without access to parallel code-switched references. This weak supervision constrains the model’s
 1238 ability to learn context-specific switching behavior beyond what is imposed by local segment similarity
 1239 and predefined constraints.

1242 Language and Domain Generalization Our study focuses on five Kenyan languages in a broadcast
1243 news context. While this setting ensures clean and aligned data, the model may not generalize to
1244 informal, multi-party, or highly emotional speech domains without further tuning or retraining.
1245

1246 Segment-Level Constraints Without Syntax Awareness Although segment replacement is guided by
1247 semantic and prosodic alignment, the model does not enforce syntactic compatibility between the
1248 inserted segment and surrounding context. This may occasionally result in grammatically awkward
1249 utterances, particularly in morphologically rich languages.
1250

1251 Speaker Identity Harmonization Is Post Hoc While a refinement step is used to harmonize speaker
1252 identity, it is applied after generation and not jointly optimized with the diffusion process. As a result,
1253 subtle speaker inconsistencies may persist across segments in certain cases.
1254

1255 Metrics May Not Capture Cultural or Pragmatic Fit Automated evaluation metrics (e.g., COMET,
1256 LaBSE) and even human Likert ratings may overlook deeper cultural or conversational appro-
1257 priateness of switches. For instance, switching at discourse boundaries or for emphasis may be
1258 underrepresented in synthetic data.
1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295