Advancing Zero-shot Text-to-Speech Intelligibility across **Diverse Domains via Preference Alignment**

Anonymous ACL submission

Abstract

Modern zero-shot text-to-speech (TTS) systems, despite using extensive pre-training, often underperform in challenging scenarios such as tongue twisters, repeated words, codeswitching, and cross-lingual synthesis, leading to intelligibility issues. This paper proposes to use preference alignment to address these challenges. Our approach leverages a newly proposed Intelligibility Preference Speech Dataset (INTP) and applies Direct Preference Optimization (DPO), along with our designed extensions, for diverse TTS architectures. After INTP alignment, in addition to intelligibility, we ob-013 serve overall improvements including naturalness, similarity, and audio quality for multiple TTS models across diverse domains. Based on 017 that, we also verify the weak-to-strong generalization ability of INTP for more intelli-018 019 gible models such as CosyVoice 2 and Ints. Moreover, we showcase the potential for further improvements through iterative alignment based on Ints. Audio samples are available at https://intalign.github.io/.

1 Introduction

011

024

037

041

Despite leveraging large-scale pre-training (Anastassiou et al., 2024; Wang et al., 2025a; Du et al., 2024b), modern zero-shot TTS systems still lack robustness during real-world applications (Sahoo et al., 2024; Neekhara et al., 2024). These systems struggle to meet even the most fundamental requirement of speech synthesis – *intelligibility* (Tan, 2023) in several scenarios, including: (1) the target text is hard to pronounce, such as tongue twisters or continuously repeated words (Neekhara et al., 2024; Anastassiou et al., 2024), which is referred to as *articulatory* cases in this paper, (2) *code-switching* cases, where the target text contains a mixture of multiple languages, and (3) cross*lingual* cases, where the languages of the target text and the reference speech differ. In these domains, existing zero-shot TTS models frequently

exhibit "hallucination" issues, such as content insertion, omission, and mispronunciation (Neekhara et al., 2024; Wang et al., 2023).

042

043

044

045

047

049

051

054

056

057

060

061

062

063

064

065

066

067

068

069

070

071

072

074

075

076

077

078

081

We attribute these intelligibility challenges primarily to the problem of out-of-distribution (OOD). For example, in cross-lingual cases, there exists a huge mismatch between monolingual pre-training and cross-lingual inference. While including such scenarios in pre-training data would be a natural solution, collecting high-quality data for challenging cases like cross-lingual synthesis remains difficult.

Motivated by the above, we propose to use preference alignment (PA) (Ouyang et al., 2022; Bai et al., 2022) to mitigate the OOD issues, and thus enhance zero-shot TTS intelligibility. The potential of this approach lies in two aspects. First, PA's customized post-training on human expected distribution can effectively mitigate the OOD issue (Ouyang et al., 2022; Xiong et al., 2024). Second, unlike TTS pre-training that requires highquality supervised data, PA needs only paired samples with relative preferences - notably, even synthetic data can lead to large improvements (Dubey et al., 2024; Yang et al., 2024b), thus significantly simplifying data collection for challenging scenarios like cross-lingual cases.

Centered on this direction, this study investigates three research problems:

- **P1**: How can we construct a high-quality intelligibility preference dataset? What prompts and base models should be selected, and how can we establish human-aligned preference pairs?
- **P2**: Unlike textual LLMs with predominantly autoregressive (AR) design, zero-shot TTS models employ diverse architectures, including ARbased (Borsos et al., 2023a; Anastassiou et al., 2024; Du et al., 2024b), Flow-Matching (FM) based (Le et al., 2023; Eskimez et al., 2024; Chen et al., 2024c), and Masked Generative Model (MGM) based (Ju et al., 2024; Wang

179

180

181

133

134

135

0

100 101

> 105 106 107

103

104

108 109

110 111 112

113 114

115 116

117

118

119

120

et al., 2025a). How can we design alignment algorithms for various architectures?

• **P3**: How well does a constructed preference dataset exhibit weak-to-strong generalization (Burns et al., 2024), i.e., the effectiveness when training more powerful models?

In this paper, we address the aforementioned problems with the following key contributions:

 \rightarrow **P1**: We establish a synthetic <u>Intelligibility</u> Preference Speech Dataset (INTP), comprising about 250K preference pairs (over 2K hours) of diverse domains. Specifically, INTP covers multiple scenarios, utilizing various TTS models for data creation. Besides, we employ several strategies to construct preference pairs, aiming to mitigate the risk of reward hacking for simple patterns (Skalse et al., 2022; Weng, 2024). Particularly, we leverage human knowledge and DeepSeek-V3 (DeepSeek-AI et al., 2024) to introduce perturbations into TTS systems, creating human-guided negative samples. In addition, when using Word Error Rate (WER) to determine intelligibility preferences, we not only consider self-comparison within a single model as in previous studies (Tian et al., 2024; Yao et al., 2025; Hussain et al., 2025), but also introduce comparisons across different models to leverage their complementary capabilities.

→ P2: We adopt the idea of Direct Preference Optimization (DPO) (Rafailov et al., 2023) to enhance various zero-shot TTS architectures. We employ the vanilla DPO algorithm for AR-based TTS models, while proposing extended versions of it for FM-based and MGM-based models. Our experiments on INTP shows that these algorithms effectively improve the intelligibility, naturalness, and overall quality of multiple state-of-the-art TTS systems, including ARS (AR-based) (Wang et al., 2025a), F5-TTS (FM-based) (Chen et al., 2024c), and MaskGCT (MGM-based) (Wang et al., 2025a).

 \rightarrow **P3**: To investigate INTP's weak-to-strong 121 generalization capability (Burns et al., 2024) on 122 more powerful base models, we research its align-123 ment effects on CosyVoice 2 (Du et al., 2024b) 124 and Ints (Appendix D). Both models are initialized 125 from textual LLMs (CosyVoice 2: from Qwen2.5, 126 0.5B (Yang et al., 2024a). Ints: from Phi-3.5-mini-128 instruct, 3.8B (Abdin et al., 2024)) and achieve superior intelligibility performance (Table 4). Our 129 experimental results verify that INTP remains ef-130 fective for these more capable models, improving 131 both intelligibility and naturalness. Additionally, 132

we showcase how to establish an *iterative* preference alignment "flywheel" of data and model improvements (Bai et al., 2022; Dubey et al., 2024; Xiong et al., 2024) based on Ints.

We will open-source all resources used in this study, including: (1) the proposed INTP and DPObased alignment codebase for various TTS models, (2) all the INTP-enhanced models based on Ints, CosyVoice 2, ARS, F5-TTS, and MaskGCT, and (3) our newly constructed zero-shot TTS evaluation sets across diverse domains.

2 INTP: Intelligibility Preference Dataset

To enhance the TTS intelligibility, this study opts for constructing a preference dataset to align (Tian et al., 2024; Yao et al., 2025; Hussain et al., 2025) rather than directly optimizing single metrics or rules such as WER (Anastassiou et al., 2024; Du et al., 2024b). This choice is motivated by two key considerations. First, through the construction of a preference dataset, we can inject human knowledge and feedback beyond WER, such as creating human-guided negative samples (Section 2.3). Second, in addition to the existing approach of constructing preference pairs from multiple samples of a single model (Tian et al., 2024; Yao et al., 2025; Hussain et al., 2025), we can leverage comparisons across different models to create preference pairs, thereby utilizing the complementary capabilities of various models (Figure 1b). These different strategies help increase diversity in the dataset, mitigating the risk of "reward hacking" that often results from the simple patterns inherent in single metrics or rules (Bai et al., 2022; Skalse et al., 2022; Weng, 2024).

Formally, we aim to construct an intelligibility preference dataset $\mathcal{D} = \{(x, y^w, y^l)\}$, where each triplet comprises a prompt x (consisting of target text x^{text} and reference speech x^{speech} for zeroshot TTS models), along with a pair of synthesized speech samples (y^w, y^l) . Here, y^w and y^l represent the preferred (positive) and dispreferred (negative) outputs conditioned on x, respectively. Statistics of the proposed INTP are presented in Table 1.

2.1 Prompt Construction

To establish a high-quality preference dataset, we aim to make the distribution of prompt x cover a wide range of domains. For the target text x^{text} , from the linguistic perspective, we design three distinct categories: (1) **Regular text**, which repre-

-							Text Type	Example
	Regular	Repeated	Code-Switching	Pronunciation- perturbed	Punctuation- perturbed	#Total	Regular	A panda eats shoots and leaves.
ARS (Wang et al., 2025a) F5-TTS (Chen et al., 2024c)	8,219	8,852 8,555	8,300 7,976	7,325 7,909	8,036 6,667	40,732	Repeated	A panda panda eats shoots and leaves and leaves.
MaskGCT (Wang et al., 2025a)	9,055	10,263	8,289	7,604	7,686	42,897	Code-Switching	熊猫吃 shoots 和 leaves。
Intra Pairs Inter Pairs	25,699 27,008	27,670 27,676	24,565 24,651	22,838 25,045	22,389 23,970	123,161 128,350	Pronunciation- perturbed	A pan duh eights shots n leafs.
#Total	52,707	55,346	49,216	47,883	46,359	251,511	Punctuation-	A panda eats, shoots, and leaves.

(a) Distribution of preference pairs, where pronunciation-perturbed and punctuationperturbed texts are introduced to create the human-guided negative samples. (b) Examples of different types for a text, *"A panda eats shoots and leaves"*.

Table 1: Intelligibility Preference dataset (INTP). There are about 250K pairs (over 2K hours) in INTP, covering various texts and speechs, multiple models, and diverse preference pairs.

182 sents the general cases for TTS systems, aimed at enhancing model intelligibility in common scenar-183 ios; (2) **Repeated text**, which contains repeated or 184 185 redundant words and phrases, specifically designed to improve TTS performance in articulatory cases; 186 and (3) Code-switching text, which incorporates 187 a mixture of different languages, intended to en-188 hance TTS capabilities in multilingual scenarios. From the semantic perspective, we collect text con-190 tent across diverse topics and domains to enrich 191 the distribution of x^{text} . For the reference speech x^{speech} , we aim to cover a wide range of speakers, 193 speaking styles, and acoustic environments. Re-194 garding the pairing of x^{text} and x^{speech} , we further 195 consider their language alignment by constructing 196 both monolingual and cross-lingual combinations (more statistics in Appendix B.1). 198

> We construct these prompt data based on the Emilia-Large (He et al., 2024, 2025), which contains real-world speech data and textual transcriptions across diverse topics, scenarios, and speaker styles. We perform stratified sampling on Emilia-Large's speech and text data to obtain multilingual prompts. We employ DeepSeek-V3 (DeepSeek-AI et al., 2024) to preprocess the sampled text, including typo correction, and use it as regular text. Based on these regular texts, we further utilize DeepSeek-V3 to transform them into different text types (as shown in Table 1b). Construction details are provided in Appendix B.1.

2.2 Model Selection

199

206

210

211

212

We utilize multiple zero-shot TTS models with diverse architectures for data synthesis to enhance INTP's diversity and generalization. Specifically, we select the following three models: (1) **ARS** (AR-based): Introduced as an autoregressive baseline by Wang et al. (2025a). and referred to as "AR + SoundStorm" in the original paper (Wang et al., 2025a). It adopts a cascaded architecture, including the autoregressive *text-to-codec* and the non-autoregressive *codec-to-waveform* (Borsos et al., 2023b). (2) **F5-TTS** (FM-based): It follows E2 TTS (Eskimez et al., 2024) and uses a flow-matching transformer (Le et al., 2023; Lipman et al., 2023) to convert the text to acoustic features directly (Chen et al., 2024c). (3) **MaskGCT** (MGM-based): Similar to ARS, MaskGCT employs a two-stage architecture. The key distinction lies in its use of an MGM in the text-to-codec stage (Wang et al., 2025a). 221

222

223

224

225

226

227

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

249

250

251

252

253

254

255

256

257

258

All the three are pre-trained on Emilia (He et al., 2024) (about 100K hours of multilingual data) and represent state-of-the-art zero-shot TTS systems across different architectures. We utilize their officially released pre-trained models (see Appendix B.2 for details) to generate data for INTP.

2.3 Preference Pairs Construction

In constructing intelligibility preference pairs, we design three categories of pairs (Figure 1):

Intra Pair These pairs are generated through model self-comparison (Figure 1a), following an approach similar to previous studies (Tian et al., 2024; Yao et al., 2025; Hussain et al., 2025). For a given prompt x, we conduct multiple samplings using the same model. Subsequently, we calculate the WER for each generation and designate the samples with the lowest and highest WER as y^w and y^l , respectively. To enlarge the gap between y^w and y^l , we employ diverse sampling hyperparameters across multiple generations from the same model. Additionally, we use a specific WER threshold to filter out pairs with insufficient performance gaps (more details in Appendix B.3.1).

Inter Pair These pairs are constructed by comparing outputs across different models (Figure 1b). The efficacy of this approach lies in leveraging the complementary strengths of various models. For



(c) Perturbed Pair Figure 1: Three kinds of preference pairs in INTP.

example, by comparing intra-pairs from different models for the same prompt, we can identify the "best of the best" samples, thereby enhancing the overall quality of positive samples in our dataset. Similar to intra pair, we also employ WER to identify intelligibility preferences for inter pairs (see Appendix B.3.2 for details).

259 260

261

263

265

272

273

277

279

281

294

298

Notably, the proposed inter-pair construction pipeline enables comparative evaluation of intelligibility performance across different models. Using this pipeline, we compared four state-of-the-art models in the field: ARS (Wang et al., 2025a), F5-TTS (Chen et al., 2024c), MaskGCT (Wang et al., 2025a), and CosyVoice 2 (Du et al., 2024b). We constructed 10K inter-pairs and analyzed the win rates of these models, as shown in Table 2. Interestingly, even ARS, the model with the lowest win rate, achieves a 4.1% success rate against the strongest model, CosyVoice 2. This finding validates our assumption regarding the complementary capabilities among various models.

Perturbed Pair In addition to the aforementioned two types of pairs which are established based on WER, we leverage human knowledge and the intelligence of DeepSeek-V3 (DeepSeek-AI et al., 2024) to create human-guided negative samples, termed perturbed pairs (Figure 1c). The main idea involves deliberately perturbing the input prompt, thereby inducing the model to generate low-quality samples (Majumder et al., 2024; Fu et al., 2024).

Specifically, we design two types of perturbation for the target text in the prompt (as shown in Table 1b): (1) **Pronunciation perturbation**: we replace certain characters of the text with easily mispronounceable alternatives. For example, given the text "*A panda eats shoots and leaves*", we can create the perturbed text "*A pan duh eights shots n leafs*". (2) **Punctuation perturbation**: we modify the punctuation, such as commas, to alter pause

	ARS	F5-TTS	MaskGCT	CosyVoice 2	Win Rate (†)
ARS	/	6.7%	7.4%	4.1%	18.3%
F5-TTS	10.4%	/	8.8%	5.9%	25.1%
MaskGCT	10.4%	8.0%	/	5.9%	24.3%
CosyVoice 2	11.9%	10.2%	10.3%	/	32.3%

* The percentage in each cell represents the proportion of cases where the model on the horizontal axis outperforms the model on the vertical axis.

* The Win Rate is calculated as the sum of values from columns 2 through 5.

Table 2: TTS Intelligibility Arena: We employ the interpair construction from INTP to compare intelligibility among four state-of-the-art zero-shot TTS models.

	ARS	F5-TTS	MaskGCT	CosyVoice 2
Positive Samples	73.0%	88.1%	90.9%	100.0%
Negative Samples	45.7%	15.8%	47.1%	75.0%
All	59.7%	53.7%	64.3%	90.4%

Table 3: Human-annotated reading accuracy (\uparrow) for four state-of-the-art zero-shot TTS models on regular texts. We use the intra-pair pipeline of INTP to generate the positive and negative samples.

299

300

301

302

303

304

305

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

329

330

patterns and prosody in the text. For example, by adding commas to the text "*A panda eats shoots and leaves*", we obtain "*A panda eats, shoots, and leaves*", where the words "*shoots*" and "*leaves*" transform from nouns in the original text to verbs, creating a significant semantic shift. The detailed process for constructing these perturbed texts is provided in Appendix B.3.3.

2.4 Human Perception Verification

After constructing INTP, we further conducted subjective evaluation to verify its alignment with human perception. For intelligibility alignment, we design a reading accuracy listening task (see Appendix F.3 for details): given a text and a speech, subjects perform binary classification to determine whether the speech accurately reads the text without any content insertion, omission, or mispronunciation. Using four state-of-the-art zero-shot TTS models, we generate 300 intra-pairs on INTP regular texts. The results in Table 3 demonstrate that INTP's preference identification for intra pairs aligns well with human judgments of intelligibility. Furthermore, comparing Tables 2 and 3 reveals that INTP's inter-pair comparisons of intelligibility across different models also effectively align with human values.

In addition to intelligibility, we also investigated how well INTP aligns with human preferences for *naturalness*, which is one of the most generalpurpose metrics for TTS (Tan, 2023). The experimental results demonstrate that the naturalness gap between positive and negative samples of INTP is

380

381

382

383

384

385

386

388

389

390

391

392

393

394

395

396

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

substantial and perceptible to human listeners. Wediscuss this finding in details in Appendix B.4.

3 Preference Alignment for Diverse Zero-Shot TTS models

In this section, we present methods for achieving preference alignment across a range of TTS models, including autoregressive based, flow-matching based, and masked generative model based architectures. Building on the framework of Direct Preference Optimization (DPO) (Rafailov et al., 2023), initially developed for AR-based models, we adapt and extend its principles to FM-based and MGMbased models.

3.1 DPO for AR Models

333

334

338

342

354

361

364

366

367

371

373

375

The main idea of reinforcement learning (RL) for preference alignment is to introduce a reward model r(x, y) to guide the model for improvement (see e.g., (Li et al., 2024)). Here y represents the output (i.e., the generated speech in zeroshot TTS), and x means the input prompt (i.e., the reference speech and the target text in zero-shot TTS). A widely adopted reward model design is based on Bradley-Terry (BT) model, which defines the probability of preferred sample y^w over dispreferred sample y^l given x as $p_{\text{BT}}(y^w \succ y^l \mid x) = \sigma(r(x, y^w) - r(x, y^l))$. We can train the reward model $r_{\phi}(x, y)$ by minimizing the negative log-likelihood of observed comparisons from the preference dataset \mathcal{D} :

$$\mathcal{L}_{\mathbf{R}} = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(r_{\phi}(x, y_w) - r_{\phi}(x, y_l) \right) \right].$$
(1)

With the given reward model, the RL optimization objective is to guide the model to maximize the expected reward while minimizing the KL-divergence from a reference distribution:

$$\max_{p_{\theta}} \mathbb{E}_{x, y \sim p_{\theta}(y|x)}[r(x, y)] - \beta D_{\mathrm{KL}}[p_{\theta}(y|x) \parallel p_{\mathrm{ref}}(y|x)],$$
(2)

where the hyperparameter β controls the strength of the regularization. As highlighted in Rafailov et al. (2023), the optimization problem in Equation 2 admits a closed form solution. This implies a direct relationship between the reward function and the policy. Substituting the reward expression into Equation 1 leads the DPO loss:

$$\mathcal{L}_{\text{DPO}} = -\mathbb{E}_{\mathcal{D}} \left[\log \sigma \left(\beta \left(\log \frac{p_{\theta}(y_w|x)}{p_{\text{ref}}(y_w|x)} - \log \frac{p_{\theta}(y_l|x)}{p_{\text{ref}}(y_l|x)} \right) \right) \right].$$
(3)

DPO enables direct preference alignment for ARbased TTS models, eliminating the need for explicit reward modeling or RL optimization. In the following subsections, we will introduce its extensions for FM-based and MGM-based TTS models.

3.2 DPO for Flow-Matching Models

The vanilla DPO algorithm is tailored for AR models, while Wallace et al. (2024) extends it to diffusion models. In this subsection, we introduce the DPO algorithm for flow-matching models, specifically demonstrating its application to optimal transport flow-matching (OT-FM), a common approach in FM-based TTS models (Le et al., 2023; Eskimez et al., 2024; Chen et al., 2024c). Given the continuous representation y of a speech sample and its corresponding condition x, OT-FM constructs a linear interpolation path between Gaussian noise $y_0 \sim \mathcal{N}(0, I)$ and the target data $y_1 = y$. Specifically, the interpolation follows $y_t = (1-t)y_0 + ty_1$, where $t \in [0, 1]$, which naturally induces a velocity field $v_{\theta}(y_t, t, x)$ that captures the constant directional derivative $\frac{dy_t}{dt} = y_1 - y_0$. OT-FM aims to learn the velocity field to match the true derivative. The corresponding loss function is defined as

$$\mathcal{L}_{\text{OT-FM}} = \mathbb{E}_{y_0, y_1, x, t} \| v_{\theta}(y_t, t, x) - (y_1 - y_0) \|_2^2, \quad (4)$$

where t is the time step that is sampled from the uniform distribution $\mathcal{U}(0, 1)$.

Inspired by Wallace et al. (2024), we rewrite the RL objective for flow-matching models. Let $p_{\theta}(y_1|y_t, t, x)$ denote our policy that predicts the target sample y_1 given the noised observation y_t at time t and condition x. We initialize from a reference flow-matching policy p_{ref} . The RL objective can be written as:

$$\max_{p_{\theta}} \mathbb{E}_{y_1 \sim p_{\theta}(y_1|x), t, x}[r(y_1, x)] - \beta \mathbb{D}_{\mathrm{KL}}[p_{\theta}(y_1|y_t, t, x) \| p_{\mathrm{ref}}(y_1|y_t, t, x)].$$
(5)

Following a similar derivation process as in DPO (we provide more details in Appendix C.2), we can obtain the loss function for flow-matching DPO:

$$\mathcal{L}_{\text{DPO-FM}} = -\mathbb{E}_{(y_1^w, y_1^l, x) \sim \mathcal{D}, t}$$
$$\log \sigma \left(\beta \left(\log \frac{p_{\theta}(y_1^w | y_t^w, t, x)}{p_{\text{ref}}(y_1^w | y_t^w, t, x)} - \log \frac{p_{\theta}(y_1^l | y_t^l, t, x)}{p_{\text{ref}}(y_1^l | y_t^l, t, x)} \right) \right), \quad (6)$$

where y_1^w and y_1^l represent the preferred and dispreferred samples from the preference dataset, respectively, while y_t^w and y_t^l are the interpolations at time t between y_1^w and y_1^l and the randomly sampled y_0^w and y_0^l . The loss can be transformed into the velocity space:

Model	ŀ	Regular	cases	Ar	ticulator	y cases	Code	e-switch	ing cases	Cro	ss-lingu	al cases		Avg	
	WER	SIM	N-CMOS	WER	SIM	N-CMOS	WER	SIM	N-CMOS	WER	SIM	N-CMOS	WER	SIM	N-CMOS
ARS	3.96	0.717	-	20.03	0.693	-	54.15	0.693	-	19.76	0.630	-	24.47	0.683	-
w/ INTP	2.32	0.727	$0.47_{\pm 0.22}$	12.83	0.713	$0.64_{\pm 0.31}$	36.91	0.698	$0.63_{\pm 0.34}$	9.57	0.632	$0.82_{\ \pm 0.28}$	15.41	0.692	$0.64_{\pm 0.12}$
F5-TTS	3.44	0.670	-	16.84	0.635	-	33.99	0.609	-	16.86	0.546	-	17.78	0.615	-
w/ INTP	2.38	0.652	$0.38 \scriptscriptstyle \pm 0.26$	12.97	0.628	$0.30 _{\pm 0.23}$	15.98	0.576	$0.67 \scriptscriptstyle \pm 0.36$	7.13	0.509	$0.47_{\ \pm 0.30}$	9.62	0.591	$0.44 \scriptstyle \pm 0.12$
MaskGCT	2.34	0.738	-	12.43	0.714	-	29.06	0.696	-	12.34	0.629	-	14.04	0.694	-
w/ INTP	2.23	0.737	$0.23_{\ \pm 0.20}$	9.13	0.722	$0.57_{\ \pm 0.36}$	19.70	0.704	$0.19_{\ \pm 0.16}$	7.87	0.633	$0.29_{\ \pm 0.18}$	9.73	0.699	$0.32_{\ \pm 0.15}$
CosyVoice 2	2.09	0.709	-	8.12	0.696	-	33.36	0.672	-	8.78	0.600	-	13.09	0.669	-
w/ INTP	1.65	0.709	$0.24_{\ \pm 0.25}$	6.87	0.696	$0.20_{\ \pm 0.16}$	28.31	0.671	$0.63_{\ \pm 0.30}$	5.39	0.603	$0.28_{\ \pm 0.31}$	10.56	0.670	$0.33_{\ \pm 0.12}$
Ints	3.14	0.688	-	12.08	0.666	-	22.88	0.646	-	9.78	0.572	-	11.97	0.643	-
w/ INTP	2.36	0.686	$0.20_{\ \pm 0.36}$	9.38	0.664	$0.11_{\pm 0.22}$	13.80	0.642	$0.20_{\ \pm 0.38}$	6.28	0.571	$0.18_{\ \pm 0.23}$	7.96	0.641	$0.17_{\pm 0.15}$

Table 4: Improvements of DPO with INTP for different models (**AR-based**: ARS (Wang et al., 2025a), CosyVoice 2 (Du et al., 2024a), and Ints (Appendix D). **FM-based**: F5-TTS (Chen et al., 2024c). **MGM-based**: MaskGCT (Wang et al., 2025a)) on diverse domains.

$$\begin{aligned} \mathcal{L}_{\text{DPO-FM}} &= -\mathbb{E}_{(y_1^w, y_1^l, x) \sim \mathcal{D}, t} \log \sigma \Big(-\beta \\ \Big(\left\| v_\theta(y_t^w, t, x) - (y_1^w - y_0^w) \right\|_2^2 - \left\| v_{\text{ref}}(y_t^w, t, x) - (y_1^w - y_0^w) \right\|_2^2 \Big) \\ &- \Big(\left\| v_\theta(y_t^l, t, x) - (y_1^l - y_0^l) \right\|_2^2 - \left\| v_{\text{ref}}(y_t^l, t, x) - (y_1^l - y_0^l) \right\|_2^2 \Big) \Big). \end{aligned}$$

This proposed algorithm can be applied to a wide range of FM-based and diffusion-based TTS models (Le et al., 2023; Eskimez et al., 2024; Shen et al., 2024). In this study, we use it to optimize F5-TTS (Chen et al., 2024c) as a representative.

3.3 DPO for Masked Generative Models

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

Masked generative model (MGM) is a type of Non-AR generative model, which is also widely adopted in speech generation, as seen in models such as NaturalSpeech 3 (Ju et al., 2024), and MaskGCT (Wang et al., 2025a). MGM aims to recover a discrete sequence $y = [z_1, z_2, ..., z_n]$ from its partially masked version $y_t = y \odot m_t$, where $m_t \in \{0, 1\}^n$ is a binary mask sampled via a schedule $\gamma(t) \in (0, 1]$. MGM is trained to predict masked tokens from unmasked tokens and condition x, modeled as $p_{\theta}(y_0 \mid y_t, x)$, optimizing the sum of the marginal cross-entropy for each unmasked token:

$$\mathcal{L}_{\text{mask}} = -\mathbb{E}_{y,x,t,m_t} \sum_{i=1}^n m_{t,i} \cdot \log p_{\theta}(z_i \mid y_t, x).$$
(8)

Using a similar derivation as in Section 3.2, we extend DPO for MGM. Let $p_{ref}(y_0 | y_t, x)$ represent the reference policy. The DPO loss for MGM is given by:

$$\mathcal{L}_{\text{DPO-MGM}} = -\mathbb{E}_{(y^w, y^l, x) \sim \mathcal{D}, t}$$
$$\log \sigma \left(\beta \left(\log \frac{p_{\theta}(y^w_0 | y^w_t, x)}{p_{\text{ref}}(y^w_0 | y^w_t, x)} - \log \frac{p_{\theta}(y^l_0 | y^l_t, x)}{p_{\text{ref}}(y^l_0 | y^l_t, x)} \right) \right).$$
(9)

Here, y_t^w and y_t^l are masked versions of y_0^w and y_0^l . Note that $p_{\theta}(y_0|y_t, x)$ corresponds to the sum of the log-probabilities of the unmasked tokens in the context of MGM. We provide more details about the derivation in Appendix C.3. In this study, we select MaskGCT (Wang et al., 2025a) as a representative to apply this proposed algorithm for its text-to-codec stage. 448

449

450

451

452

453

4 Experiments

Evaluation Data We evaluate zero-shot TTS sys-454 tems across diverse domains in both English and 455 Chinese languages. Based on SeedTTS's evalua-456 tion samples (Anastassiou et al., 2024) (which are 457 widely used and also serve as the evaluation set for 458 the pre-trained models of ARS (Wang et al., 2025a), 459 F5-TTS (Chen et al., 2024c), MaskGCT (Wang 460 et al., 2025a), and CosyVoice 2 (Du et al., 2024b) 461 in this study), we construct evaluation sets across 462 four distinct domains: (1) Regular cases: We use 463 SeedTTS test-en (1,000 samples) and SeedTTS 464 test-zh datasets (2,000 samples). (2) Articula-465 tory cases: These involve tongue twisters and re-466 peated texts. For Chinese, we use SeedTTS test-467 hard, while for English, we use reference speech 468 prompts of SeedTTS test-en, and employ Deepseek-469 V3 (DeepSeek-AI et al., 2024) to construct the ar-470 ticulatory texts like SeedTTS test-hard. There are 471 800 samples in total. (3) Code-switching cases: 472 These target texts are a mixture of English and Chi-473 nese. Based on SeedTTS test-en and test-zh, we 474 keep their reference speech prompts unchanged, 475 and adopt Deepseek-V3 to transform their texts 476 into code-switching style. There are 1,000 samples 477 in total. (4) Cross-lingual cases: We construct two 478 types of cross-lingual samples: zh2en (500 sam-479 ples) and *en2zh* (500 samples). The zh2en means 480 Chinese reference speech (from SeedTTS test-zh) 481 with English target text (from SeedTTS test-en). 482 Similarly for en2zh. The detailed distribution of 483



Figure 2: Subjective evaluation of intelligibility and speaker similarity for models before and after INTP alignment.

these sets is presented in Table 9, Appendix F.1.

484

485

486

487

488

489

491

492

493

494

495

496

497

498

499

505

506

507

508

510

511

512

513

515

516

517

519

521

523

Evaluation Metrics For objective metrics, we evaluate the intelligibility (WER, \downarrow), speaker similarity (SIM, \uparrow), and overall speech quality (UT-MOS (Saeki et al., 2022), \uparrow). Specifically, for WER, we employ Whisper-large-v3 (Radford et al., 2023) for English, and Paraformer-zh (Gao et al., 2022, 2023) for Chinese and code-switching texts. For SIM, we compute the cosine similarity between the WavLM TDNN (Chen et al., 2022) speaker embeddings of generated samples and the reference speeches. For subjective metrics, we employ Comparative Mean Opinion Score (rated from -2 to 2) to evaluate naturalness (N-CMOS, \uparrow), use reading accuracy (Section 2.4) to evaluate intelligibility, and use A/B Testing to compare speaker similarity between the generated samples before and after intelligibility alignment. Detailed descriptions of all the metrics are provided in Appendix F.

4.1 Effect of DPO with INTP

To verify the effectiveness of DPO with INTP for existing TTS models, we conduct alignment experiments with multiple models. In addition to ARS, F5-TTS, and MaskGCT, which were used in constructing the INTP dataset, we also introduce two more powerful models in terms of intelligibility: CosyVoice 2 (Du et al., 2024b) and Ints (Appendix D), to validate INTP's weak-to-strong generalization capability. The experimental results are presented in Table 4, including results on the objective WER, SIM, and the subjective naturalness CMOS.

We observe three key findings from Table 4: (1) Across different evaluation cases, while almost all models demonstrate strong intelligibility performance in regular cases (WER < 4.0), they struggle significantly with articulatory, code-switching, and cross-lingual cases. We show some hallucinated outputs for these domains on our demo website. (2) Comparing across models, the proposed Ints achieves the best average intelligibility performance across all cases (WER of 11.97), highlighting the strength of using a textual LLM as the initialization of large-scale TTS model (Du et al., 2024b). (3) Through DPO with INTP, all models, including the more intelligible CosyVoice 2 and Ints that are out of the INTP distribution, show improvements in both intelligibility (WER) and naturalness (N-CMOS), and display comparable performance for speaker similarity (SIM). 524

525

526

527

528

529

530

531

532

533

534

535

536

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

555

556

557

558

559

560

561

563

Furthermore, we randomly sample 300 samples for subjective evaluation, including assessments of reading accuracy and A/B testing of speaker similarity before and after INTP alignment (see Appendix F.3 for details). The results in Figure 2 demonstrate that INTP alignment enhances all five models in terms of both intelligibility (higher reading accuracy in Figure 2a) and speaker similarity (more Tie/Win percentages in Figure 2b).

4.2 Effect of Different Data within INTP

To investigate the impact of different distributions within INTP, we conduct ablation studies from multiple perspectives. In Table 5, we present three groups of experiments on ARS: the effect of data across different text types, across different models, and the effect of different negative samples. Additional results, including the effect of data across different languages are provided in Appendix G.

We observe three key findings from Table 5: (1) Group 1 demonstrates that different scenarios require customized post-training data. For instance, repeated data proves particularly effective for articulatory cases, while pronunciation-perturbed data significantly improves pronunciation accuracy and WER in cross-lingual cases (see our demo website for details). Moreover, utilizing data from multiple scenarios (i.e., the complete INTP) yields the best overall improvements. (2) Group 2 reveals that model improvement can be achieved through alignment using synthetic data, regardless of whether

Model	R	Regular o	cases	Art	iculator	y cases	Code	-switchi	ng cases	Cro	ss-lingu	al cases		Avg	
model	WER	SIM	UTMOS	WER	SIM	UTMOS	WER	SIM	UTMOS	WER	SIM	UTMOS	WER	SIM	UTMOS
				Group 1	: Effect	of Data acro	oss Diffe	rent Tex	t Types						
ARS (Wang et al., 2025a)	3.96	0.717	3.145	20.03	0.693	2.915	54.15	0.693	3.045	19.76	0.630	3.120	24.47	0.683	3.056
w/ Regular	2.45	0.727	3.200	17.41	0.706	3.000	37.52	0.701	3.110	9.66	0.638	3.200	16.76	0.693	3.128
w/ Repeated	2.33	0.725	3.225	12.88	0.711	3.050	39.74	0.701	3.150	10.96	0.636	3.235	16.48	0.693	3.165
w/ Code-switching	2.32	0.729	3.220	17.67	0.704	3.050	34.20	0.695	3.140	8.69	0.633	3.215	15.72	0.690	3.156
w/ Pronunciation-perturbed	2.21	0.720	3.250	17.76	0.693	3.075	35.99	0.687	3.185	8.24	0.617	3.285	16.05	0.679	3.199
w/ Punctuation-perturbed	2.46	0.722	3.240	17.35	0.699	3.020	42.73	0.694	3.160	10.94	0.624	3.255	18.37	0.684	3.169
w/ INTP	2.32	0.727	3.210	12.83	0.713	3.035	36.91	0.698	3.145	9.57	0.632	3.250	15.41	0.692	3.160
				Group	2: Effec	et of Data ac	ross Diff	ferent M	odels						
ARS (Wang et al., 2025a)	3.96	0.717	3.145	20.03	0.693	2.915	54.15	0.693	3.045	19.76	0.630	3.120	24.47	0.683	3.056
w/ ARS pairs	2.56	0.717	3.200	13.05	0.705	3.015	40.91	0.691	3.125	11.07	0.622	3.225	16.90	0.684	3.141
w/ MaskGCT pairs	2.37	0.724	3.210	16.85	0.700	3.010	37.41	0.692	3.105	8.83	0.625	3.200	16.37	0.685	3.131
w/ F5-TTS pairs	2.46	0.721	3.210	14.99	0.705	3.035	38.77	0.690	3.115	10.01	0.621	3.225	16.56	0.684	3.146
w/ Intra pairs	2.33	0.721	3.200	15.29	0.705	3.015	37.99	0.687	3.115	9.36	0.624	3.200	16.24	0.684	3.133
w/ Inter pairs	2.25	0.726	3.180	15.42	0.703	2.965	38.69	0.697	3.065	10.61	0.631	3.170	16.74	0.689	3.095
w/ INTP	2.32	0.727	3.210	12.83	0.713	3.035	36.91	0.698	3.145	9.57	0.632	3.250	15.41	0.692	3.160
				Group	o 3: Effe	ct of Differe	ent Nega	tive Sam	ples						
ARS (Wang et al., 2025a)	3.96	0.717	3.145	20.03	0.693	2.915	54.15	0.693	3.045	19.76	0.630	3.120	24.47	0.683	3.056
w/ Regular (SFT)*	3.28	0.716	3.165	20.03	0.685	2.935	48.73	0.691	3.065	17.25	0.630	3.165	22.32	0.680	3.083
w/ Regular*	2.45	0.727	3.200	17.41	0.706	3.000	37.52	0.701	3.110	9.66	0.638	3.200	16.76	0.693	3.128
w/ Pronunciation-perturbed*	2.21	0.720	3.250	17.76	0.693	3.075	35.99	0.687	3.185	8.24	0.617	3.285	16.05	0.679	3.199
w/ Punctuation-perturbed*	2.46	0.722	3.240	17.35	0.699	3.020	42.73	0.694	3.160	10.94	0.624	3.255	18.37	0.684	3.169

* The positive samples in these four experiments are identical. w/ Regular (SFT) refers to supervised fine-tuning using positive samples only, excluding negative samples. w/ Regular employs WER-based negative samples, while the other two utilize our proposed human-guided negative samples.

Table 5: Effect of different data within INTP for ARS.

Model	Preference Data		Regular	cases	Art	iculator	y cases	Code	-switchi	ng cases	Cro	ss-lingu	al cases		Avg	
		WER	SIM	UTMOS	WER	SIM	UTMOS	WER	SIM	UTMOS	WER	SIM	UTMOS	WER	SIM	UTMOS
Ints	-	3.14	0.688	3.175	12.08	0.666	3.025	22.88	0.646	3.045	9.78	0.572	3.150	11.97	0.643	3.099
Ints v1	INTP	2.36	0.686	3.205	9.38	0.664	3.060	13.80	0.642	3.125	6.28	0.571	3.230	7.96	0.641	3.155
Ints v2	Ints v1 generated	2.21	0.686	3.210	8.48	0.660	3.085	12.33	0.643	3.140	5.40	0.567	3.250	7.10	0.639	3.171

Table 6: Iterative Preference Alignment for Ints.

it's generated by the model itself or other models.
Besides, the intra-pairs and inter-pairs are complementary for model improvements. (3) Group 3 shows that using only positive samples from INTP for supervised fine-tuning (SFT) can already improve quality. Building upon this, incorporating negative samples for preference learning leads to even more substantial gains.

4.3 Iterative Intelligibility Alignment

564

565

569

571

574

575

576

579

583

584

586

587

Furthermore, we explore how to establish an *iterative* preference alignment, i.e., data and model flywheel (Bai et al., 2022; Dubey et al., 2024; Xiong et al., 2024). We investigate two rounds of alignment based on Ints, where Ints v1 (INTP-aligned model) is used to generate new preference data for training Ints v2, following a similar *cadence* of data collection as (Bai et al., 2022). To prepare Ints v1 generated preference data, we sample a challenging prompt subset from INTP and adopt the same pipeline as INTP to construct preference pairs (see Appendix D.2 for details). The results of this iterative alignment are shown in Table 6. We can observe that compared to Ints v1, Ints v2 yields additional improvements across all scenarios, which demonstrates that effectiveness of iterative alignment. However, we observe that the magnitude of improvement in the second round is notably smaller than the first round. We suspect this indicates that the upper bound of iterative alignment is largely determined by the base model's inherent capabilities, suggesting future research should focus on base models with higher potential.

5 Conclusion

In this work, we focus on the intelligibility issues of modern zero-shot TTS systems across diverse domains, especially in hard-to-pronounce texts, codeswitching, and cross-lingual synthesis. We propose to address these challenges using preference alignment with our newly constructed INTP dataset, which contains diverse preference pairs determined through model self-comparison, cross-model comparison, and human guidance. We employ DPO and design special extensions to significantly improve various TTS architectures, while demonstrating INTP's weak-to-strong generalization capability and establishing an iterative preference alignment flywheel with more powerful base models.

605

606

607

609

610

588

589

590

611 Limitations

While our approach demonstrates significant improvements in zero-shot TTS intelligibility across 613 diverse domains, several limitations remain. Al-614 though INTP covers multiple challenging scenar-615 ios, it may not fully capture all edge cases, such 616 as specialized jargon or rare language pairs. Fu-617 ture work could expand to more low-resource lan-618 guages and niche domains. Besides, constructing INTP and conduct alignment experiments on large models like Ints require substantial computational 621 resources, potentially limiting accessibility.

Potential Risks

624

625

628

630

631

634 635

641

642

644

645

647

651

652

653

654

656

657

662

The proposed method introduces several risks that warrant consideration. Enhanced TTS systems could be exploited to generate deceptive content (e.g., deepfake audio), posing ethical challenges. Robust safeguards and watermarking mechanisms are critical for deployment. While INTP uses public datasets, real-world applications may risk incorporating sensitive or copyrighted speech data, requiring strict governance protocols.

References

Marah I Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat S. Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Caio César Teodoro Mendes, Weizhu Chen, Vishrav Chaudhary, Parul Chopra, Allie Del Giorno, Gustavo de Rosa, Matthew Dixon, Ronen Eldan, Dan Iter, Amit Garg, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis, Dongwoo Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Chen Liang, Weishung Liu, Eric Lin, Zeqi Lin, Piyush Madan, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Xia Song, Masahiro Tanaka, Xin Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Michael Wyatt, Can Xu, Jiahang Xu, Sonali Yadav, Fan Yang, Ziyi Yang, Donghan Yu, Chengruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. Phi-3 technical report: A highly capable language model locally on your phone. arXiv preprint, abs/2404.14219.

Chen, Zhuo Chen, Ziyi Chen, Jian Cong, Lelai Deng, Chuang Ding, Lu Gao, Mingqing Gong, Peisong Huang, Qingqing Huang, Zhiying Huang, Yuanyuan Huo, Dongya Jia, Chumin Li, Feiya Li, Hui Li, Jiaxin Li, Xiaoyang Li, Xingxing Li, Lin Liu, Shouda Liu, Sichao Liu, Xudong Liu, Yuchen Liu, Zhengxi Liu, Lu Lu, Junjie Pan, Xin Wang, Yuping Wang, Yuxuan Wang, Zhen Wei, Jian Wu, Chao Yao, Yifeng Yang, Yuanhao Yi, Junteng Zhang, Qidi Zhang, Shuo Zhang, Wenjie Zhang, Yang Zhang, Zilin Zhao, Dejian Zhong, and Xiaobin Zhuang. 2024. Seed-tts: A family of high-quality versatile speech generation models. *arXiv preprint*, abs/2406.02430.

664

665

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

- Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M Tyers, and Gregor Weber. 2019. Common voice: A massivelymultilingual speech corpus. *arXiv preprint arXiv:1912.06670*.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Chris Olah, Benjamin Mann, and Jared Kaplan. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint*, abs/2204.05862.
- Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matthew Sharifi, Dominik Roblek, Olivier Teboul, David Grangier, Marco Tagliasacchi, and Neil Zeghidour. 2023a. Audiolm: A language modeling approach to audio generation. *IEEE ACM Trans. Audio Speech Lang. Process.*, 31:2523–2533.
- Zalán Borsos, Matthew Sharifi, Damien Vincent, Eugene Kharitonov, Neil Zeghidour, and Marco Tagliasacchi. 2023b. Soundstorm: Efficient parallel audio generation. *arXiv preprint*, abs/2305.09636.
- Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner, Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, Ilya Sutskever, and Jeffrey Wu. 2024. Weakto-strong generalization: Eliciting strong capabilities with weak supervision. In *ICML*. OpenReview.net.
- Chen Chen, Yuchen Hu, Wen Wu, Helin Wang, Eng Siong Chng, and Chao Zhang. 2024a. Enhancing zero-shot text-to-speech synthesis with human feedback. *arXiv preprint*, abs/2406.00654.
- Jingyi Chen, Ju-Seung Byun, Micha Elsner, and Andrew Perrault. 2024b. Dlpo: Diffusion model loss-guided reinforcement learning for fine-tuning text-to-speech diffusion models. *arXiv preprint*, abs/2405.14632.
- Philip Anastassiou, Jiawei Chen, Jitong Chen, Yuanzhe

Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al. 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.

721

722

724

728

729

730

731

732

733

735

737

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

761

763

764

766

767

770

771

772

773 774

775

776

777

778

779

781

- Yushen Chen, Zhikang Niu, Ziyang Ma, Keqi Deng, Chunhui Wang, Jian Zhao, Kai Yu, and Xie Chen.
 2024c. F5-TTS: A fairytaler that fakes fluent and faithful speech with flow matching. *arXiv preprint*, abs/2410.06885.
- Geoffrey Cideron, Sertan Girgin, Mauro Verzetti, Damien Vincent, Matej Kastelic, Zalán Borsos, Brian McWilliams, Victor Ungureanu, Olivier Bachem, Olivier Pietquin, Matthieu Geist, Léonard Hussenot, Neil Zeghidour, and Andrea Agostinelli. 2024. Musicrl: Aligning music generation to human preferences. In *ICML*. OpenReview.net.
- Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. Flashattention: Fast and memory-efficient exact attention with io-awareness. Advances in Neural Information Processing Systems, 35:16344–16359.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi,

Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. 2025. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning.

782

783

784

785

786

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, and Wangding Zeng. 2024. Deepseek-v3 technical report. arXiv preprint, abs/2412.19437.
- Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi. 2023. High fidelity neural audio compression. *Trans. Mach. Learn. Res.*, 2023.
- Zhihao Du, Qian Chen, Shiliang Zhang, Kai Hu, Heng Lu, Yexin Yang, Hangrui Hu, Siqi Zheng, Yue Gu, Ziyang Ma, Zhifu Gao, and Zhijie Yan. 2024a. Cosyvoice: A scalable multilingual zero-shot textto-speech synthesizer based on supervised semantic tokens. *arXiv preprint*, abs/2407.05407.
- Zhihao Du, Yuxuan Wang, Qian Chen, Xian Shi, Xiang Lv, Tianyu Zhao, Zhifu Gao, Yexin Yang, Changfeng Gao, Hui Wang, Fan Yu, Huadai Liu, Zhengyan Sheng, Yue Gu, Chong Deng, Wen Wang, Shiliang Zhang, Zhijie Yan, and Jingren Zhou. 2024b. Cosyvoice 2: Scalable streaming speech synthesis with large language models. *arXiv preprint*, abs/2412.10117.

- 841 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, 842 Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-852 lonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, 862 Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, 866 Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate 871 Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. 2024. The llama 3 herd of models. arXiv preprint, abs/2407.21783. 873
 - Sefik Emre Eskimez, Xiaofei Wang, Manthan Thakker, Canrun Li, Chung-Hsien Tsai, Zhen Xiao, Hemin Yang, Zirun Zhu, Min Tang, Xu Tan, Yanqing Liu, Sheng Zhao, and Naoyuki Kanda. 2024. E2 TTS: embarrassingly easy fully non-autoregressive zeroshot TTS. In *SLT*. IEEE.

876

877

878

879

884

894

896

900

- Deqing Fu, Tong Xiao, Rui Wang, Wang Zhu, Pengchuan Zhang, Guan Pang, Robin Jia, and Lawrence Chen. 2024. TLDR: token-level detective reward model for large vision language models. *arXiv preprint*, abs/2410.04734.
- Xiaoxue Gao, Chen Zhang, Yiming Chen, Huayun Zhang, and Nancy F. Chen. 2024. Emo-dpo: Controllable emotional speech synthesis through direct preference optimization. *arXiv preprint*, abs/2409.10157.
- Zhifu Gao, Zerui Li, Jiaming Wang, Haoneng Luo, Xian Shi, Mengzhe Chen, Yabin Li, Lingyun Zuo, Zhihao Du, Zhangyu Xiao, et al. 2023. Funasr: A fundamental end-to-end speech recognition toolkit. arXiv preprint arXiv:2305.11013.
- Zhifu Gao, Shiliang Zhang, Ian McLoughlin, and Zhijie Yan. 2022. Paraformer: Fast and accurate parallel transformer for non-autoregressive end-to-end speech recognition. *arXiv preprint arXiv:2206.08317*.
- Hao-Han Guo, Kun Liu, Fei-Yu Shen, Yi-Chen Wu, Feng-Long Xie, Kun Xie, and Kai-Tuo Xu. 2024. Fireredtts: A foundation text-to-speech framework for

industry-level generative speech applications. *arXiv* preprint, abs/2409.03283.

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

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

- Tingwei Guo, Cheng Wen, Dongwei Jiang, Ne Luo, Ruixiong Zhang, Shuaijiang Zhao, Wubo Li, Cheng Gong, Wei Zou, Kun Han, et al. 2021. Didispeech: A large scale mandarin speech corpus. In *ICASSP 2021-*2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 6968– 6972. IEEE.
- Haorui He, Zengqiang Shang, Chaoren Wang, Xuyuan Li, Yicheng Gu, Hua Hua, Liwei Liu, Chen Yang, Jiaqi Li, Peiyang Shi, Yuancheng Wang, Kai Chen, Pengyuan Zhang, and Zhizheng Wu. 2024. Emilia: An extensive, multilingual, and diverse speech dataset for large-scale speech generation. In *SLT*. IEEE.
- Haorui He, Zengqiang Shang, Chaoren Wang, Xuyuan Li, Yicheng Gu, Hua Hua, Liwei Liu, Chen Yang, Jiaqi Li, Peiyang Shi, et al. 2025. Emilia: A large-scale, extensive, multilingual, and diverse dataset for speech generation. *arXiv preprint*, 2501.15907.
- Yuchen Hu, Chen Chen, Siyin Wang, Eng Siong Chng, and Chao Zhang. 2024. Robust zero-shot text-tospeech synthesis with reverse inference optimization. *arXiv preprint*, abs/2407.02243.
- Shehzeen Hussain, Paarth Neekhara, Xuesong Yang, Edresson Casanova, Subhankar Ghosh, Mikyas T. Desta, Roy Fejgin, Rafael Valle, and Jason Li. 2025. Koel-tts: Enhancing llm based speech generation with preference alignment and classifier free guidance. *arXiv preprint*, abs/2502.05236.
- Zeqian Ju, Yuancheng Wang, Kai Shen, Xu Tan, Detai Xin, Dongchao Yang, Eric Liu, Yichong Leng, Kaitao Song, Siliang Tang, et al. 2024. Naturalspeech 3: Zero-shot speech synthesis with factorized codec and diffusion models. In *Forty-first International Conference on Machine Learning*.
- Jacob Kahn, Morgane Riviere, Weiyi Zheng, Evgeny Kharitonov, Qiantong Xu, Pierre-Emmanuel Mazaré, Julien Karadayi, Vitaliy Liptchinsky, Ronan Collobert, Christian Fuegen, et al. 2020. Libri-light: A benchmark for asr with limited or no supervision. In ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 7669–7673. IEEE.
- Matthew Le, Apoorv Vyas, Bowen Shi, Brian Karrer, Leda Sari, Rashel Moritz, Mary Williamson, Vimal Manohar, Yossi Adi, Jay Mahadeokar, and Wei-Ning Hsu. 2023. Voicebox: Text-guided multilingual universal speech generation at scale. In *NeurIPS*.
- Ziniu Li, Tian Xu, Yushun Zhang, Zhihang Lin, Yang Yu, Ruoyu Sun, and Zhi-Quan Luo. 2024. Remax: A simple, effective, and efficient reinforcement learning method for aligning large language models. In *Fortyfirst International Conference on Machine Learning*.

- 956 957 960 961 962 963 964 965 966 967 969 970 971 972 973 974 975 977 978 979 982 985 987 991 992 997
- 998
- 999 1000 1001
- 1002
- 1004 1005 1006

1009

1010 1011 1012

- Huan Liao, Haonan Han, Kai Yang, Tianjiao Du, Rui Yang, Qinmei Xu, Zunnan Xu, Jingquan Liu, Jiasheng Lu, and Xiu Li. 2024. BATON: aligning text-to-audio model using human preference feedback. In IJCAI, pages 4542–4550. ijcai.org.
- Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. 2023. Flow matching for generative modeling. In ICLR. Open-Review.net.
- Navonil Majumder, Chia-Yu Hung, Deepanway Ghosal, Wei-Ning Hsu, Rada Mihalcea, and Soujanya Poria. 2024. Tango 2: Aligning diffusion-based text-toaudio generations through direct preference optimization. In ACM Multimedia, pages 564-572. ACM.
- Paarth Neekhara, Shehzeen Hussain, Subhankar Ghosh, Jason Li, Rafael Valle, Rohan Badlani, and Boris Ginsburg. 2024. Improving robustness of llm-based speech synthesis by learning monotonic alignment. In INTERSPEECH. ISCA.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In NeurIPS.
- 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. IEEE.
- Puyuan Peng, Po-Yao Huang, Shang-Wen Li, Abdelrahman Mohamed, and David Harwath. 2024. Voicecraft: Zero-shot speech editing and text-to-speech in the wild. arXiv preprint arXiv:2403.16973.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In International conference on machine learning, pages 28492-28518. PMLR.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. In NeurIPS.
- Takaaki Saeki, Detai Xin, Wataru Nakata, Tomoki Koriyama, Shinnosuke Takamichi, and Hiroshi Saruwatari. 2022. UTMOS: utokyo-sarulab system for voicemos challenge 2022. In INTERSPEECH, pages 4521-4525. ISCA.
- Pranab Sahoo, Prabhash Meharia, Akash Ghosh, Sriparna Saha, Vinija Jain, and Aman Chadha. 2024. A comprehensive survey of hallucination in large language, image, video and audio foundation models. In EMNLP (Findings), pages 11709–11724. Association for Computational Linguistics.

Kai Shen, Zeqian Ju, Xu Tan, Eric Liu, Yichong Leng, Lei He, Tao Qin, Sheng Zhao, and Jiang Bian. 2024. Naturalspeech 2: Latent diffusion models are natural and zero-shot speech and singing synthesizers. In ICLR. OpenReview.net.

1013

1014

1015

1017

1018

1019

1020

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1041

1042

1047

1048

1049

1050

1051

1052

1053

1054

1056

1057

1058

1060

1061

1062

1063

1064

1065

1066

- Joar Skalse, Nikolaus Howe, Dmitrii Krasheninnikov, and David Krueger. 2022. Defining and characterizing reward gaming. In NeurIPS, volume 35, pages 9460-9471.
- Anonymous Preprint Submission. 2025. Dualcodec: A low-frame-rate, semantically-enhanced neural audio codec for speech generation. Available on OpenReview
- Xu Tan. 2023. Neural Text-to-Speech Synthesis. Springer.
- Jinchuan Tian, Chunlei Zhang, Jiatong Shi, Hao Zhang, Jianwei Yu, Shinji Watanabe, and Dong Yu. 2024. Preference alignment improves language modelbased TTS. arXiv preprint, abs/2409.12403.
- Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam, Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. 2024. Diffusion model alignment using direct preference optimization. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 8228-8238.
- 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. *arXiv preprint*, abs/2301.02111.
- Yuancheng Wang, Haoyue Zhan, Liwei Liu, Ruihong Zeng, Haotian Guo, Jiachen Zheng, Qiang Zhang, Xueyao Zhang, Shunsi Zhang, and Zhizheng Wu. 2025a. Maskgct: Zero-shot text-to-speech with masked generative codec transformer. In ICLR. OpenReview.net.
- Yuancheng Wang, Jiachen Zheng, Junan Zhang, Xueyao Zhang, Huan Liao, and Zhizheng Wu. 2025b. Metis: A foundation speech generation model with masked generative pre-training. arXiv preprint arXiv:2502.03128.
- Lilian Weng. 2024. Reward hacking in reinforcement learning. lilianweng.github.io.
- Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. 2024. Iterative preference learning from human feedback: Bridging theory and practice for RLHF under klconstraint. In ICML. OpenReview.net.
- Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. 2023. Imagereward: Learning and evaluating human preferences for text-to-image generation. In NeurIPS.

- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, 1068 Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan 1069 Li, Dayiheng Liu, Fei Huang, Guanting Dong, Hao-1070 ran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin 1072 Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, 1077 Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, 1078 Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, 1079 Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing 1080 Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, 1083 Zhifang Guo, and Zhihao Fan. 2024a. Qwen2 technical report. arXiv preprint, abs/2407.10671.
 - An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2024b. Qwen2.5 technical report. *arXiv preprint*, abs/2412.15115.

1088

1089 1090

1091

1092

1093

1096

1098

1100

1102

1103

1104

1105

1106

1107 1108

1109

1110

1111

1112

1113

1114

- Jixun Yao, Yuguang Yang, Yu Pan, Yuan Feng, Ziqian Ning, Jianhao Ye, Hongbin Zhou, and Lei Xie. 2025. Fine-grained preference optimization improves zeroshot text-to-speech. *arXiv preprint*, abs/2502.02950.
 - Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi. 2021. Soundstream: An end-to-end neural audio codec. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 30:495–507.
 - Dong Zhang, Zhaowei Li, Shimin Li, Xin Zhang, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. 2024. Speechalign: Aligning speech generation to human preferences. In *NeurIPS*.
- Xueyao Zhang, Xiaohui Zhang, Kainan Peng, Zhenyu Tang, Vimal Manohar, Yingru Liu, Jeff Hwang, Dangna Li, Yuhao Wang, Julian Chan, Yuan Huang, Zhizheng Wu, and Mingbo Ma. 2025. Vevo: Controllable zero-shot voice imitation with self-supervised disentanglement. In *ICLR*. OpenReview.net.

A Related Work

1116

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

Zero-Shot Text to Speech Given a target text 1117 and a reference speech as input, zero-shot TTS 1118 systems aim to synthesize the target text while 1119 mimicking the reference style. Modern zero-shot 1120 TTS systems include AR approaches (Wang et al., 1121 1122 2023; Peng et al., 2024; Anastassiou et al., 2024; Guo et al., 2024; Du et al., 2024a,b; Zhang et al., 1123 2025) that model discrete speech tokens (Zeghidour 1124 et al., 2021; Défossez et al., 2023), and Non-AR 1125 approaches that either model continuous represen-1126 tations using diffusion (Shen et al., 2024) or flow 1127 matching (Le et al., 2023; Eskimez et al., 2024; 1128 Chen et al., 2024c), or model discrete tokens using 1129 masked generative models (Borsos et al., 2023b; 1130 Ju et al., 2024; Wang et al., 2025a,b). While these 1131 systems, trained on large-scale datasets (He et al., 1132 2024; Kahn et al., 2020; He et al., 2025), show 1133 excellent intelligibility in regular cases (Anastas-1134 siou et al., 2024; Panayotov et al., 2015; Du et al., 1135 2024b), they still struggle with intelligibility in 1136 real-world scenarios. 1137

Alignment for Speech Generation Alignment via post-training has demonstrated its effectiveness in the generation of text (Ouyang et al., 2022; Bai et al., 2022), vision (Xu et al., 2023; Fu et al., 2024), speech (Zhang et al., 2024; Anastassiou et al., 2024; Du et al., 2024b), music (Cideron et al., 2024), and sound effects (Majumder et al., 2024; Liao et al., 2024). In speech generation, existing works have employed preference alignment to enhance multiple aspects of speech, including intelligibility (Anastassiou et al., 2024; Du et al., 2024b; Tian et al., 2024), speaker similarity (Anastassiou et al., 2024; Du et al., 2024b; Tian et al., 2024), emotion controllability (Anastassiou et al., 2024; Gao et al., 2024), and overall quality (Zhang et al., 2024; Chen et al., 2024a; Hu et al., 2024; Chen et al., 2024b; Yao et al., 2025; Hussain et al., 2025). For intelligibility, previous studies choose WER as the optimization objective, either directly employing it as a reward model (Anastassiou et al., 2024; Du et al., 2024b) or centering around it to construct preference pairs (Tian et al., 2024; Yao et al., 2025; Hussain et al., 2025).

1161However, the existing research exhibits two main1162limitations. First, in constructing intelligibility1163preference dataset, current works rely solely on1164a single model to generate data (Tian et al., 2024;1165Yao et al., 2025; Hussain et al., 2025), neglecting1166comparisons across different models. Additionally,

beyond the objective WER, the potential of lever-1167 aging human knowledge or feedback to construct 1168 preference pairs remains unexplored. Second, most 1169 existing work has focused primarily on optimiz-1170 ing AR-based (Zhang et al., 2024; Anastassiou 1171 et al., 2024; Du et al., 2024b; Tian et al., 2024) or 1172 diffusion-based (Chen et al., 2024b) TTS models, 1173 leaving open the question of how to design effec-1174 tive alignment algorithms for other architectural 1175 paradigms, such as FM-based and MGM-based 1176 TTS models. 1177

B Construction Details of INTP

B.1 Prompt Construction

We construct English and Chinese prompt data, both based on the Emilia-Large dataset (He et al., 2024, 2025), which contains diverse real-world speech data across various topics, recording scenarios, and speaking styles.

Reference Speech We perform stratified sampling on Emilia-Large's speech data based on its metadata such as topics and tags to cover diverse acoustic conditions. Considering the memory constraints of existing zero-shot TTS models during inference, we only select samples with durations not exceeding 12 seconds.

Target Text Similarly to reference speech, we perform stratified sampling based on Emilia-Large's metadata to cover diverse semantic topics. We select speech samples with durations between 5 and 22 seconds, and use their corresponding textual transcriptions as the target text data source.

We utilize DeepSeek V3 (DeepSeek-AI et al., 2025) to preprocess the sampled textual transcriptions, such as typo correction and punctuation mark normalization, and use the processed text as regular text in INTP. Specifically, we use the following instruction for DeepSeek V3 to conduct text preprocessing:

System Prompt:

I obtained a text from an audio file based on some ASR models. Please help me clean it up (e.g., correct typos, add proper punctuation marks, and make the sentences semantically coherent). Note: (1) You can modify, add, or replace words that better fit the context to ensure semantic coherence. (2) Please only return the cleaned-up result without any explanation.

User Prompt (Example): a panda eats shoes and leaves

System Output (Example): A panda eats shoots and leaves.

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

Furthermore, we employ DeepSeek V3 to transform the regular text into different types. To generate Chinese-English-mixed code-switching texts:

System Prompt:

请你把这句话,转换成一个中文、英文混合的 codeswitching 版本。注意:你只需要返回给我转换后的 结果,不需要任何解释。

User Prompt (Example): A panda eats shoots and leaves.

System Output (Example): 熊猫吃 shoots 和 leaves。

1210

1211

1212

1213

1214

1215

1216

To generate punctuation-perturbed texts:

System Prompt:

假设你是一个 Text To Speech (TTS)领域的专家, 现在,让我们对一个 TTS 系统进行攻击。具体地: 我输入一个文本,请你修改这条文本里面的若干词 语,从而使 TTS 系统更容易出错。例如:你可以修 改为把某些字修改为容易读错的形近字、把多音字 做替换,等等,但你不要增加和删除原有的文本。 注意:你只需要返回给我转换后的结果,不需要任 何解释。

例子1: 【我的输入】我今天很高兴 【你的输出】窝锦添狠搞醒

例子2

【我的输入】目前,爱心人士正在种作寄养的小猫 已经五个月大了。而本人的种作寄养申请单需要进 一步审核。为了避免小猫多次转手,治疗者们对小 猫的种作寄养提出了严格要求:申请人需年满二十 三岁。 【你的输出】幕前,爱信人士正在重作寄扬的削猫 已经伍个月大了。而本人的重作寄扬神情但需要进 一步审核。为了闭面削猫多次转售,治理者们对削 猫的重作寄扬提出了阉割要求:申情人需年慢贰拾

叁岁。 例子3:

【我的输入】 And the idea of standing all by himself in a crowded market, to be pushed and hired by some big, strange farmer, was very disagreeable. Why not sing that high note and grow potatoes?

【你的输出】And the eye dear of standing awl bye himself in a crowd dead market, two bee pushed and high red buy sum big, strange far mer, was vary dis agreeable. Y knot sing that hi note and grow poe eight toes?

User Prompt (Example): A panda eats shoots and leaves.

System Output (Example): A pan duh eights shots n leafs.

To generate repeated text and punctuationperturbed text, we leverage DeepSeek V3 to create executable Python scripts that implement rulebased word repetition and random punctuation modification. These scripts will be included in our future open-source repository.

Combination between Speech and Text Based 1218 on the language of reference speech and target 1219 text data, we design four balanced combination 1220 categories: monolingual combinations (en2en and 1221 *zh2zh*) and cross-lingual combinations (*zh2en* and 1222 en2zh), where zh2en denotes Chinese reference 1223 speech with English target text, and similarly for 1224 others. For each text type shown in Table 1a (Reg-1225 ular, Repeated, Code-Switching, Pronunciation-1226 perturbed, and Punctuation-perturbed), we con-1227 struct 12K prompts. 1228

1217

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1256

1257

1258

1259

1260

1261

B.2 Model Selection

- **ARS** (Wang et al., 2025a): We use the original checkpoint (pre-trained on Emilia) provided by the authors.
- **F5-TTS** (Chen et al., 2024c): We use the officially released checkpoint¹ for INTP data generation.
- MaskGCT (Wang et al., 2025a): We use the officially released checkpoint² for INTP data generation.

In addition to these three models used for INTP construction, we also investigate INTP's effectiveness on **CosyVoice 2** and **Ints**. For CosyVoice 2, we conduct alignment experiments using its officially released checkpoint³ as the base model. Details of the pre-trained models of Ints are provided in Appendix D.

B.3 Preference Pairs Construction

B.3.1 Intra Pair

For each model and prompt, we perform five samplings and construct intra pairs based on their WER comparisons. To maximize the performance gap between positive and negative samples, we employ two strategies. First, we use diverse hyperparameters during the five generations to increase sample diversity, selecting the generation with the lowest WER as positive samples and the highest WER as negative samples. Second, we apply a threshold to filter out pairs where the WER gap between positive and negative samples is less than 6.0.

Specifically, for ARS's five samplings, we set top k to 20 and top p to 1.0, while using different temperatures of 0.4, 0.6, 0.8, 1.0, and 1.2. For F5-TTS

¹https://huggingface.co/SWivid/F5-TTS

²https://huggingface.co/amphion/MaskGCT

³https://github.com/FunAudioLLM/CosyVoice

A+2	A+1	Tie	B+1	B+2
10.9%	29.0%	15.0%	32.4%	12.6%
* For each raters indicat better to moder Tie indice	ch pair, we in random tes that sam than B, whi ately better licates no p	present the order, labe ple A's natu le A+1 ind than B, sir erceptible of	two sample eled as A a uralness is s licates that s nilar for B- difference.	es to human and B. A+2 significantly sample A is +2 and B+1.

Table 7: Human naturalness preference for 1,000 pairs from INTP regular text domain.

	Naturalness	Naturalness	Naturalness
	Winner	Tie	Loser
INTP winner	72%	15%	13%

Table 8: Agreement between INTP preference and human naturalness preference.

and MaskGCT, we use the generated speech target duration as the sampling hyperparameter. Denoting the "ground truth" duration⁴ as *d*, we employ five different duration parameters: 0.8*d*, 0.9*d*, 1.0*d*, 1.1*d*, and 1.2*d*.

B.3.2 Inter Pair

1263

1264

1265

1266

1269

1270

1271

1272

1273

1275

1276

1277

1278

1280

1281

1282

1283

1284

1285

1286 1287

1288

1289

We construct inter pairs based on the intra pairs established in Appendix B.3.1. For a given prompt, we denote model A's intra pair as (y_A^w, y_A^l) and model B's intra pair as (y_B^w, y_B^l) . We construct inter pairs through three types of comparisons: between y_A^w and y_B^w , between y_A^w and y_B^l , and between y_A^l and y_B^w . Note that we exclude comparisons between y_A^l and y_B^w . Note that we exclude comparisons between y_A^l and y_B^w . Note that we exclude comparisons between y_A^l and y_B^l to ensure high quality of positive samples. We apply the same WER threshold as in Appendix B.3.1 to filter out pairs with small performance gaps.

B.3.3 Perturbed Pair

The instructions used to prompt DeepSeek V3 for obtaining pronunciation-perturbed and punctuationperturbed texts are shown in Appendix B.1. Specifically, we only use data from INTP's regular text domain to construct perturbed pairs.

B.4 Human Verification

In Section 2.4, we evaluated INTP's alignment with human intelligibility perception. In this section, we investigate the alignment between INTP and human naturalness preferences. Specifically, we design a **naturalness preference** annotation task (Ap-1290 pendix F.3). We randomly sample 1,000 pairs from 1291 INTP's regular text domain for human annotation, 1292 with results shown in Table 7 and 8. The results 1293 reveal two key findings: First, 85% of INTP pairs 1294 exhibit distinguishable naturalness preferences (Tie 1295 rate of 15% in Table 7). Additionally, INTP's pref-1296 erence determination shows strong agreement with 1297 human naturalness preferences (72% agreement 1298 rate between INTP winners and naturalness win-1299 ners in Table 8). These results suggest that INTP 1300 can also serve as a foundation dataset for natural-1301 ness preference alignment in future research. 1302

C Details of the Derivation

C.1 DPO for AR Models

Starting from Equation 2, Rafailov et al. (2023) demonstrate that the optimization problem admits a closed-form solution. Specifically, the optimal policy $p_{\theta}^*(y|x)$ that maximizes the RL objective is given by:

1303

1304

1305

1306

1307

1308

1310

1311

1312

1313

1315

1316

1317

1318

1320

1321

1322

$$p_{\theta}^{*}(y|x) = \frac{1}{Z(x)} p_{\text{ref}}(y|x) \exp\left(\frac{1}{\beta}r(x,y)\right), \quad (10)$$

where Z(x) is the partition function ensuring normalization. This establishes a direct relationship between the reward function and the policy:

$$r(x,y) = \beta \log \frac{p_{\theta}^*(y|x)}{p_{\text{ref}}(y|x)} + \beta \log Z(x).$$
(11) 1314

Substituting this reward expression (Equation 11) into the reward modeling loss function (Equation 1) leads the DPO loss (Equation 3), which we represent here as:

$$\mathcal{L}_{\text{DPO}} = -\mathbb{E}_{\mathcal{D}}\left[\log\sigma\left(\beta\left(\log\frac{p_{\theta}(y_w|x)}{p_{\text{ref}}(y_w|x)} - \log\frac{p_{\theta}(y_l|x)}{p_{\text{ref}}(y_l|x)}\right)\right)\right].$$
1319

C.2 DPO for Flow-Matching Models

Starting from Equation 5, which we represent here as:

$$\max_{p_{\theta}} \mathbb{E}_{y_1 \sim p_{\theta}(y_1|x), t, x}[r(y_1, x)] - \beta \mathbb{D}_{\mathrm{KL}}[p_{\theta}(y_1|y_t, t, x) \| p_{\mathrm{ref}}(y_1|y_t, t, x)].$$
1323

Similar to the derivation in DPO (Rafailov et al.,13242023) and Wallace et al. (2024), we obtain the1325closed-form solution for the optimal policy as:1326

$$p_{\theta}^{*}(y_{1}|y_{t},t,x) = \frac{1}{Z(y_{t},t,x)} p_{\text{ref}}(y_{1}|y_{t},t,x) \exp\left(\frac{1}{\beta}r(y_{1},x)\right),$$
(12) 1327

⁴Since we use Emilia-Large's transcription data as target text in our prompt construction process (Appendix B.1), we refer to the original speech duration corresponding to this transcription as the "ground truth" duration.

1328 where $Z(y_t, t, x)$ is the partition function ensuring 1329 normalization. We can then express the reward 1330 model $r(y_1, x)$ as:

$$r(y_1, x) = \beta \log \frac{p_{\theta}^*(y_1|y_t, t, x)}{p_{\text{ref}}(y_1|y_t, t, x)} + \beta \log Z(y_t, t, x).$$
(13)

1332Similarly, substituting this reward expression1333(Equation 13) into the reward modeling loss func-1334tion (Equation 1) leads to the DPO loss for OT-FM1335(Equation 6), which we represent here:

1331

1336

1344

1345

1346

1347

1348

1351

$$\begin{split} \mathcal{L}_{\text{DPO-FM}} &= -\mathbb{E}_{(y_1^w, y_1^l, x) \sim \mathcal{D}, t} \\ &\log \sigma \left(\beta \left(\log \frac{p_\theta(y_1^w | y_t^w, t, x)}{p_{\text{ref}}(y_1^w | y_t^w, t, x)} - \log \frac{p_\theta(y_1^l | y_t^l, t, x)}{p_{\text{ref}}(y_1^l | y_t^l, t, x)} \right) \right). \end{split}$$

1337Reviewing the training objective of OT-FM (Equa-1338tion 4), we find that it is equivalent to fitting a1339Gaussian likelihood. In other words, the induced1340likelihood can be interpreted as:

1341
$$p_{\theta}(y_1 \mid y_t, t, x) \propto \exp\left(-\frac{1}{\beta} \|v_{\theta}(y_t, t, x) - (y_1 - y_0)\|_2^2\right)$$

similarly, for the reference policy, we have:

1343
$$p_{\rm ref}(y_1 \mid y_t, t, x) \propto \exp\left(-\frac{1}{\beta} \|v_{\rm ref}(y_t, t, x) - (y_1 - y_0)\|_2^2\right)$$

Here, β serves as an inverse temperature (or noise variance), and the normalization constants cancel out when taking the ratio. By taking the logarithm of the ratio between the learned policy and the reference policy, we obtain:

$$\log \frac{p_{\theta}(y_1 \mid y_t, t, x)}{p_{\text{ref}}(y_1 \mid y_t, t, x)} = -\frac{1}{\beta} \Big(\|v_{\theta}(y_t, t, x) - (y_1 - y_0)\|_2^2 - \|v_{\text{ref}}(y_t, t, x) - (y_1 - y_0)\|_2^2 \Big).$$

1350 Multiplying both sides by β results in:

$$\beta \log \frac{p_{\theta}(y_1 \mid y_t, t, x)}{p_{\text{ref}}(y_1 \mid y_t, t, x)} = -\left(\|v_{\theta}(y_t, t, x) - (y_1 - y_0)\|_2^2 - \|v_{\text{ref}}(y_t, t, x) - (y_1 - y_0)\|_2^2 \right)$$

1352By substituting the log-ratio formulation into Equa-1353tion 6, we can transform the DPO loss for OT-FM1354into a form related to the velocity, as shown in1355Equation 7, which is represented as:

$$\mathcal{L}_{\text{DPO-FM}} = -\mathbb{E}_{(y_1^w, y_1^l, x) \sim \mathcal{D}, t} \log \sigma \Big(-\beta \\ \Big(\| v_\theta(y_t^w, t, x) - (y_1^w - y_0^w) \|_2^2 - \| v_{\text{ref}}(y_t^w, t, x) - (y_1^w - y_0^w) \|_2^2 \Big) \\ - \Big(\| v_\theta(y_t^l, t, x) - (y_1^l - y_0^l) \|_2^2 - \| v_{\text{ref}}(y_t^l, t, x) - (y_1^l - y_0^l) \|_2^2 \Big) \Big).$$

C.3 DPO for Masked Generative Models

Similar to flow-matching, let $p_{\theta}(y_0 \mid y_t, x)$ denote the policy to be optimized, and $p_{ref}(y_0 \mid y_t, x)$ the reference policy. We can rewrite the RL objective for MGM as follows: 1361

$$\max_{p_{\theta}} \quad \mathbb{E}_{y_{0} \sim p_{\theta}(y_{0}|x), t, x} \left[r(y_{0}, x) \right] \\ -\beta \mathbb{D}_{\mathrm{KL}} \left[p_{\theta}(y_{0}|y_{t}, x) \, \| \, p_{\mathrm{ref}}(y_{0}|y_{t}, x) \right].$$
(14) 1362

1357

1363

1364

1365

1366

1368

1370

1371

1372

1374

1375

1376

1377

1378

1379

1380

We can also derive the closed-form solution for the optimal policy:

$$p_{\theta}^{*}(y_{0}|y_{t},x) = \frac{1}{Z(y_{t},x)} p_{\text{ref}}(y_{0}|y_{t},x) \exp\left(\frac{1}{\beta}r(y_{0},x)\right),$$
(15)

and express the reward model as follows:

$$r(y_0, x) = \beta \log \frac{p_{\theta}^*(y_0|y_t, x)}{p_{\text{ref}}(y_0|y_t, x)} + \beta \log Z(y_t, x), \quad (16)$$

where $Z(y_t, x)$ is the partition function ensuring normalization. Also, substituting this reward expression (Equation 16) into the reward modeling loss function (Equation 1) leads to the DPO loss for MGM:

$$\mathcal{L}_{\text{DPO-MGM}} = -\mathbb{E}_{(y^w, y^l, x) \sim \mathcal{D}, t}$$

$$\log \sigma \left(\beta \left(\log \frac{p_{\theta}(y^w_0 | y^w_t, x)}{p_{\text{ref}}(y^w_0 | y^w_t, x)} - \log \frac{p_{\theta}(y^l_0 | y^l_t, x)}{p_{\text{ref}}(y^l_0 | y^l_t, x)} \right) \right).$$
(17) 13

Here, y_t^w and y_t^l are masked versions of y_0^w and y_0^l generated via the mask schedule $\gamma(t)$. Note that $p_{\theta}(y_0|y_t, x)$ corresponds to the sum of the log-probabilities of the unmasked tokens in the context of MGM.

D Ints: Intelligibility-enhanced Speech Language Model

Ints is an intelligibility-enhanced speech language 1381 model. It follows a two-stage generation paradigm 1382 like (Anastassiou et al., 2024; Du et al., 2024a; 1383 Wang et al., 2025a): in the first stage, it uses an AR 1384 model to generate discrete speech tokens, while in 1385 the second stage, it employs a flow matching model 1386 to generate mel-spectrograms from speech tokens. We use the first-layer tokens from DualCodec (Sub-1388 mission, 2025) as the modeling target for the first 1389 stage of Ints, due to its efficient compression rep-1390 resentation (12.5Hz tokens for 24kHz speech) and 1391 rich semantic information. Particularly, the first-1392 stage AR model is directly initialized from a large 1393 language model while extending the vocabulary 1394 to include speech tokens. The codebook size of speech tokens is 16,384. Specifically, in this work, 1396

1400

1401

1402 1403

1404

1406

1407

1408

1409

1410 1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421 1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437

1438

1439

1440

1441

1405

we use the 3.8B Phi-3.5-mini-instruct⁵ (Abdin et al., 2024), motivated by scaling up model size and leveraging the rich textual semantic knowledge.

D.1 TTS Instruction Design

We format the input as a text-to-speech instruction concatenated with speech tokens. The input sequence is represented as:

$$[I, T, < | ext{startofspeech} | >, S, < | ext{endofspeech} | >]$$

where **I** is the instruction prefix (e.g., "Please speak the following text out loud"), and T and S denote the text and speech token sequences, respectively. The special tokens < |startofspeech| > and < endofspeech > mark the boundaries of the speech token sequence.

During the inference stage for zero-shot TTS, the input sequence is represented as:

 $[\boldsymbol{I}, \boldsymbol{T}_{\text{prompt}}, \boldsymbol{T}_{\text{target}}, < |\text{startofspeech}| >, \boldsymbol{S}_{\text{prompt}}]$

to generate the target speech tokens S_{target} . Here, $T_{\text{prompt}}, T_{\text{target}}, S_{\text{prompt}}$ are placeholders for the prompt text, target text, and prompt speech tokens, respectively.

D.2 Training data

We pre-train Ints on Emilia (He et al., 2024), which consists of about 100K hours of multilingual data. Following this, we use INTP alignment to obtain Ints v1. Ints v1 is then used to generate new preference data, which are employed to train Ints v2 for iterative alignment. We select prompts from the repeated and code-switching samples of INTP, which can be considered a more challenging subset of prompts. For each prompt, we use the same INTP intra-pair pipeline in Appendix B.3.1 to construct preference pairs.

Ε **Training Details**

All of our experiments are conducted on 8 NVIDIA H100 80GB-GPUs. Unless stated otherwise, we use the AdamW optimizer with $\beta_1 = 0.9, \beta_2 =$ 0.999 and train for one epoch. For each model, we provide more detailed information about the experiments:

• **ARS**: We use a learning rate of 5e - 6 with a warmup of 4,000 steps and an inverse square root learning scheduler. For DPO, we use the hyperparameter $\beta = 0.1$.

⁵https://huggingface.co/microsoft/Phi-3.5-mini-instruct

	Lang	uages	#Total
Regular	en	zh	3,000
guim	1,000	2,000	
Articulatory	en	zh	800
	400	400	
Code-switching	en2mixed	zh2mixed	1 000
code switching	500	500	
Cross-lingual	zh2en	en2zh	1.000
cross ingun	500	500	

Table 9: Statistics of the proposed evaluation sets in four scenarios (en: English, zh: Chinese, mixed: mixture of English and Chinese, zh2en: Chinese reference speech with English target text. Similarly for en2mixed, zh2mixed, and en2zh).

• **F5-TTS**: We use a learning rate of 8e - 6 with a warmup of 4,000 steps and an inverse square root learning scheduler. For DPO, we use the hyperparameter $\beta = 1,000$.

1442

1443

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

1458

1459

1460

1461

1462

1463

1464

1465

1466

1467

1468

1469

- MaskGCT: We use a learning rate of 5e-6 with a warmup of 4,000 steps and an inverse square root learning scheduler. For DPO, we use the hyperparameter $\beta = 10$.
- CosyVoice 2: We use a learning rate of 5e 6with a warmup of 4,000 steps and an inverse square root learning scheduler. For DPO, we use the hyperparameter $\beta = 0.1$.
- Ints: We use a learning rate of 5e 6 with a warmup of 4,000 steps and an inverse square root learning scheduler. For DPO, we use the hyperparameter $\beta = 0.1$. We use flash attention (Dao et al., 2022) and bfloat16 for training.

F **Evaluation Details**

F.1 Evaluation Data

Our evaluation sets are based on SeedTTS testen and SeedTTS test-zh datasets⁶. The SeedTTS test-en set includes 1,000 samples from the Common Voice dataset (Ardila et al., 2019), while the SeedTTS test-zh set comprises 2,000 samples from the DiDiSpeech dataset (Guo et al., 2021). We also provide the detailed distribution of our proposed sets in Table 9.

F.2 Objective Evaluation Metrics

For objective metrics, we evaluate the intelligibil-1470 ity (WER), speaker similarity (SIM), and overall 1471 speech quality (UTMOS (Saeki et al., 2022)): 1472

⁶https://github.com/BytedanceSpeech/seed-tts-eval

					On	English Ev	aluation	Samples	8						
Model	1	Regular	(en)	Ar	ticulator	ry (en)	Code-s	witching	g (en2mixed)	Cros	s-lingua	l (zh2en)	Avg		
	WER	SIM	UTMOS	WER	SIM	UTMOS	WER	SIM	UTMOS	WER	SIM	UTMOS	WER	SIM	UTMOS
ARS (Wang et al., 2025a)	3.55	0.682	3.560	15.98	0.675	3.400	48.59	0.629	3.190	15.22	0.697	3.150	20.84	0.671	3.325
w/ en2en	1.96	0.697	3.690	13.42	0.685	3.570	35.18	0.641	3.270	8.19	0.692	3.300	14.19	0.679	3.458
w/ zh2zh	2.76	0.692	3.660	13.90	0.687	3.550	36.65	0.644	3.260	8.92	0.694	3.320	15.06	0.679	3.448
w/ en2zh, zh2en	2.32	0.694	3.700	11.78	0.684	3.580	35.17	0.645	3.290	7.00	0.700	3.330	14.07	0.681	3.475
w/ all	2.35	0.695	3.680	13.76	0.686	3.560	33.53	0.642	3.240	7.38	0.704	3.310	14.26	0.682	3.448
					On	Chinese Ev	aluation	Samples	5						
Model	1	Regular	(<i>zh</i>)	Ar	ticulato	ry (<i>zh</i>)	Code-s	witching	g (zh2mixed)	Cros	s-lingua	l (en2zh)		Avg	
	WER	SIM	UTMOS	WER	SIM	UTMOS	WER	SIM	UTMOS	WER	SIM	UTMOS	WER	SIM	UTMOS
ARS (Wang et al., 2025a)	4.37	0.752	2.730	24.07	0.711	2.430	59.71	0.756	2.900	24.30	0.563	3.090	28.61	0.696	2.788
w/ en2en	2.68	0.761	2.760	21.68	0.727	2.530	48.84	0.757	2.990	12.48	0.566	3.140	21.42	0.703	2.855
w/ zh2zh	2.41	0.760	2.740	19.51	0.727	2.490	47.99	0.755	3.010	12.73	0.565	3.110	20.16	0.702	2.838
w/ en2zh, zh2en	2.49	0.762	2.740	22.92	0.715	2.490	41.00	0.757	3.000	11.76	0.573	3.160	19.54	0.702	2.848
w/ all	2.62	0.759	2.720	21.06	0.725	2.440	41.50	0.760	2.980	11.95	0.572	3.090	19.78	0.704	2.808

Table 10: Effect of different languages within INTP for ARS. In these experiments, we use only the **Regular** part of INTP for training.

- WER: We employ Whisper-large-v3⁷(Radford et al., 2023) for English texts, and Paraformer-zh⁸(Gao et al., 2022, 2023) for Chinese and codeswitching texts.
 - **SIM**: We compute the cosine similarity between the WavLM TDNN⁹(Chen et al., 2022) speaker embeddings of generated samples and the prompt samples.
 - **UTMOS**: We use the pretrained UTMOS strong learner following the official implementation¹⁰.

F.3 Subjective Evaluation

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1500

1501

1502

We consider four different settings: regular, articulatory, code-switching, and cross-lingual. Each setting is evaluated in two languages, resulting in 10 samples per language. This setup yields a total of 80 pairs. These 80 pairs are evaluated across 5 different systems (ARS, F5-TTS, MaskGCT, CosyVoice 2, and Ints), leading to a total of 400 pairs. We engage 20 participants in the evaluation process, ensuring that each sample is assessed at least three times.

We conduct subjective evaluations from three perspectives: intelligibility (reading accuracy), naturalness (N-CMOS), and speaker similarity (A/B Testing). We have developed an automated subjective evaluation interface, as shown in Figure 3 and Figure 4. For each item to be evaluated, users will see three components: the System Interface, the Questionnaire, and the Evaluation Criteria.

⁸https://huggingface.co/funasr/paraformer-zh

i <u>tep 1.</u> Listen Audio	Task ID #3
ත් ⁹ Ran	dom Task
Speech A	Speech B
1 2 3 4 5 6 7 8 9 10 11 12 四 有时候,我们的心就像被撕裂了一般,感受到那种深深的痛勉,仿佛 无论怎么努力都无法弥补内心的空洞。	11 2 3 4 5 5 5 7 8 9 10 11 12 12 15 7 16 17 18 19 10 11 12 15 15 15 15 15 15 15 15 15 15 15 15 15
D 00:00.00 / 00:12.76	D0:00.00 / 00:12.76
tep 2. Assessment	
Is any reading error? (insertion, nission, or mispronunciation) Speech A 🔵 Has Err	ror 🔿 No Error Speech B 🔿 Has Error 🔿 No Error
Which speech sounds more natural? A +2	A +1 Tie B +1 B +2

Figure 3: User interface for intelligibility and naturalness evaluation.

Intelligibility (Reading Accuracy):

• **System Interface:** Users listen to the speech audio and compare it to the provided target text to assess whether the speech matches the text.

1503

1504

1505

1506

1507

1508

1509

1510

1511

1512

1513

1514

1515

1516

1517

- Questionnaire: Users are asked, "Is any reading error? (insertion, omission, or mispronunciation)"
- Evaluation Criteria: The evaluation is binary: "No Error" (the speech matches the text) or "Has Error" (the speech does not match the text).

Naturalness (N-CMOS):

- **System Interface:** Users listen to two speech samples, A and B, to compare their naturalness.
- **Questionnaire:** Users are asked, "Which speech sounds more natural?"
- Evaluation Criteria: Options include A +2 (Sample A is much more natural), A +1 (Sample A is slightly more natural), Tie (Both are equally natural), B +1 (Sample B is slightly more natural).
 1521 1522
 1522

⁷https://huggingface.co/openai/whisper-large-v3

⁹https://github.com/microsoft/UniSpeech/tree/main/

downstreams/speaker_verification

¹⁰https://github.com/sarulab-speech/UTMOS22

ේ Ranc	dom Task
Reference Speaker	
↓> + ++	
□ · · · · · · · · · · · · · · · · · · ·	3
▷ 00:00	.00 / 00:04.46
Speech A	Speech B
+++ +++ ++++++++++++++++++++++++++++++	HI NO HI H- HIN H- HIN HIN HI HIN HI
1 2 3 4 5 6 7 8 9 10 11 12	1 2 3 4 5 6 7 8 9 10 11 1
百有时候,我们的心就像被撕裂了一般,感受到那种深深的痛楚,仿佛 无论怎么努力都无法弥补内心的空洞。	百 有时候,我们的心就像被撕裂了一般,感受到那种深深的痛楚,仿 无论怎么努力都无法弥补内心的空洞。
00:00.00 / 00:12.76	D0:00.00 / 00:12.76

Figure 4: User interface for speaker similarity evaluation.

Speaker Similarity (A/B Testing):

1523

1524 1525

1526

1527

1528

1529

1530

1531

1532

1533

1534

1535

- **System Interface:** Users listen to two speech samples, A and B, to evaluate their similarity to the speech of the reference speaker.
- Questionnaire: Users are asked, "Which speech sounds more like the reference speaker's style?"
- Evaluation Criteria: Options include A +2 (Sample A is much more similar), A +1 (Sample A is slightly more similar), Tie (Both are equally similar), B +1 (Sample B is slightly more similar), and B +2 (Sample B is much more similar).

G Effect of Data across Different Languages within INTP

We present the effect of different languages within 1536 INTP in Table 10. The results reveal three key find-1537 ings: (1) Data from all languages can contribute 1538 1539 to improvements across diverse domains for ARS. (2) Interestingly, using only English post-training 1540 data (w/ en2en) could also improve performance on 1541 Chinese evaluation samples, and vice versa, demon-1542 strating that the proposed alignment algorithm en-1543 hances the model's foundation capability in intel-1544 ligibility. (3) Furthermore, we again observe the 1545 effectiveness of preference alignment's customized 1546 feature: when aiming to improve performance on 1547 cross-lingual cases, directly constructing data from 1548 the cross-lingual distribution yields the most signif-1549 icant gains. 1550