

Fake it to make it: Using synthetic data to remedy the data shortage in joint multimodal speech-and-gesture synthesis

Anonymous CVPR submission

Paper ID *****

Abstract

Although humans engaged in face-to-face conversation simultaneously communicate both verbally and non-verbally, methods for joint and unified synthesis of speech audio and co-speech 3D gesture motion from text are a new and emerging field. These technologies hold great promise for more human-like, efficient, expressive, and robust synthetic communication, but are currently held back by the lack of suitably large datasets, as existing methods are trained on parallel data from all constituent modalities. Inspired by student-teacher methods, we propose a straightforward solution to the data shortage, by simply synthesising additional training material. Specifically, we use unimodal synthesis models trained on large datasets to create multimodal (but synthetic) parallel training data, and then pre-train a joint synthesis model on that material. In addition, we propose a new synthesis architecture that adds better and more controllable prosody modelling to the state-of-the-art method in the field. Our results confirm that pre-training on large amounts of synthetic data improves the quality of both the speech and the motion synthesised by the multimodal model, with the proposed architecture yielding further benefits when pre-trained on the synthetic data.

1. Introduction

Human beings are embodied, and we use a wide gamut of the expressions afforded by our bodies to communicate. In concert with the lexical and non-lexical (prosodic) components of speech, humans also leverage gestures realised by face, head, arm, finger, and body motion – all driven by a shared, underlying communicative intent [58] – to improve face-to-face communication [30, 66].

Research into automatically recreating different kinds of human communicative behaviour, whether it be speech audio from text [85], or gesture motion from speech [92], have a long history, as these are key enabling technologies for, e.g., virtual agents, game characters, and social robots

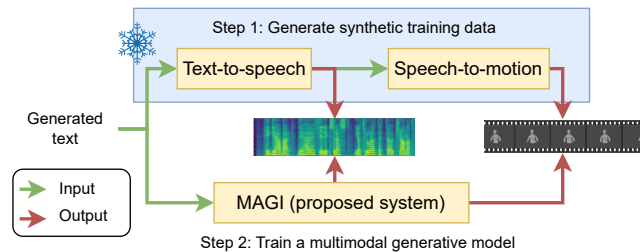


Figure 1. MAGI: Multimodal Audio and Gesture, Integrated

[14, 41, 57, 68]. The advent of deep learning has led to an explosion of research in the two fields [54, 66, 83]. Gesture synthesis, in particular, has been shown to benefit from access to both lexical and acoustic representations of speech [3, 42, 43, 104]. That said, joint and simultaneous synthesis of both speech and gesture communication (pioneered in [78]) remains severely under-explored. This despite the fact that simultaneously generating both modalities together not only better emulates how humans produce communicative expressions, but also offers a stepping stone towards creating non-redundant gestures that can complement and even replace speech, like human gestures do [34]. On top of this, recent research efforts towards integrating the synthesis of the two modalities have demonstrated improvements in coherent [6, 62], compact [62, 94], jointly and rapidly learnable [61], convincing [61, 62], and cross-modally appropriate [62] synthesis of speech and 3D gestures from text.

The current state of the art in joint multimodal speech-and-gesture synthesis, Match-TTSG [62], achieves strong performance via modern techniques such as conditional flow matching (OT-CFM) [51] with U-Net Transformer [91] encoders [77]. However, there still remains a noticeable gap between synthesised model output and recordings of natural human speech and gesticulation [62]. This contrasts with recent breakthroughs in “generative AI”, which can synthesise text [2, 13], images [77], and speech audio [80, 84] that all are nigh indistinguishable from those created by humans. The critical difference is that whereas those strong models for synthesising single modalities benefit from training on vast amounts of data (cf. [27]), exist-

ing parallel datasets of speech audio, text transcriptions, and human motion are radically smaller. This is especially true if we require good motion quality (which at present generally necessitates high-end 3D motion capture) and speech audio with a spontaneous character and quality suitable for speech synthesis. The state-of-the-art joint synthesis system demonstrated in [62] was thus trained on 4.5 hours of parallel speech and gesture data from [22]; larger parallel corpora exist [49, 53], but exhibit some quality issues (cf. [44]) and do not exceed 100 hours, a far cry from the corpora used to train leading generative AI systems. It stands to reason that multimodal synthesis systems could gain substantially from overcoming the limitations imposed by training only on presently available parallel corpora.

In this paper, we propose two improvements to the state-of-the-art multimodal speech-and-gesture synthesis:

1. We pre-train a joint speech-and-gesture synthesis model on a large parallel corpus of *synthetic* training data created using leading text, text-to-speech, and speech-to-gesture systems (Fig. 1). This provides a straightforward way to let multimodal models benefit from advances in data and systems for unimodal synthesis.
2. We extend [62] with a probabilistic duration model (similar to [48]) and individual models of pitch and energy (similar to [75]). This enables more lifelike and more controllable synthetic expression.

The resulting joint synthesis system is orders of magnitude smaller and faster than the models used for synthesising the pre-training data. Our subjective evaluations show that the proposed pre-training on synthetic data improves the speech as well as the gestures created by a joint synthesis system, and that the architectural modifications further benefit a system pre-trained on large synthetic data and also enable output control. For examples of model output, please see our anonymous webpage at cvprhumogen24.github.io/MAGI/; code will be released with future versions of the paper.

2. Background

In this section, we review synthesis of text, speech audio, and 3D gesture motion, along with existing work in multimodal speech-and-gesture synthesis. For each task, we state how the methods relate to our contributions and briefly discuss how synthetic data can improve synthesis models.

2.1. Text generation

The rise of large language models (LLMs) has brought revolutionary improvements to text generation. Transformer-based [91] LLMs using Generative Pretrained Transformers (GPTs) [71] like [2, 13, 88] are capable of generating text virtually indistinguishable from that written by humans.

The critical methodological advances for LLMs are pre-training on vast amounts of diverse data, coupled with fine-tuning on a small amount of high-quality, in-domain mate-

rial, e.g., via Reinforcement Learning from Human Feedback (RLHF) [9]. This methodology of pre-training foundation models followed by fine-tuning on the best data has been validated to give excellent results across several modalities [11, 111]. In this paper, we for the first time use that methodology in joint speech-and-gesture synthesis.

Fine-tuned LLMs allow generating of diverse text samples for many domains through *prompting* the model, i.e., providing a written text prompt at runtime describing the output to generate. Prompting has been useful for many tasks including creating synthetic dialogue datasets [1] and selecting appropriate gestures based on verbal utterances [28]. We use this ability to create an arbitrarily large material of conversational text sentences in the style of a given speaker/corpus as a basis for our synthetic-data creation.

2.2. Speech synthesis

Recent advancements in deep generative modelling have significantly improved text-to-speech (TTS) [83], achieving levels of naturalness that rival recorded human speech [80, 84]. TTS approaches are primarily divided into two broad classes: autoregressive (AR) and non-autoregressive (NAR) architectures. AR architectures produce acoustic outputs sequentially, using mechanisms such as neural cross-attention [10, 15, 50, 79, 110] or neural transducers [59, 60, 101] to connect inputs symbols to the outputs. Conversely, non-autoregressive models [25, 36, 37, 48, 63, 69, 75, 112] generate the entire utterance in parallel. The NAR approach is typically faster, especially on GPUs, but AR methods (which invest more computation into synthesis) often have the edge in synthesis quality.

Recently, there has been a trend [10, 12, 15, 46, 93] to quantise audio waveforms into discrete tokens [16, 46], and then adapt an LLM-like autoregressive approach (e.g., with GPTs) to learn to model these audio tokens on large datasets. Synthesised token sequences can subsequently be converted back to audio [81]. Speaker and style adaptation can be achieved by seeding (prompting) the model with an audio snippet, something we leverage to create diverse stochastic synthetic training data for our work.

LLM-like TTS can give exceptional results when trained on large datasets, but models risk confabulating (similar to well-known issues with LLMs) and getting trapped in feedback loops due to the autoregression [10, 15]. Our paper therefore describes a pipeline for mitigating these problems when creating synthetic training data at scale.

In NAR TTS, it has been found that conditioning the TTS on the output of a model of prosodic properties, e.g., per-phone pitch and energy, can benefit synthesis [67, 75, 112]. This furthermore affords control over speech output by replacing or manipulating the prosodic features prior to synthesis. Especially important for convincing prosody are the durations of the synthesised speech sounds. It has been

169 shown [37, 40] that probabilistic modelling of durations can
170 substantially improve deep generative TTS. This appears
171 especially useful for speech uttered spontaneously in conver-
172 sation, as considered here, due to its highly diverse and
173 non-deterministic prosodic structure [47]. Inspired by these
174 advances, we introduce a probabilistic duration model cou-
175 pled with explicit pitch and energy models into the mul-
176 timodal synthesis architecture. Better duration modelling
177 should help create speech rhythm and timings that allow
178 adequate time for gesture-preparation phases, so that beat-
179 gesture strokes can be distinct and synchronised with the
180 speech. Improved control will not only affect the output
181 speech but also the gestures we generate with it.

182 2.3. Gesture synthesis

183 Like TTS, deep learning has led to a boom in 3D ges-
184 ture synthesis from speech text and/or audio [66]. The list
185 of deep generative techniques considered includes GANs
186 [95, 96], normalising flows [4, 5], VAEs [23], VQ-VAEs
187 [102, 103], combinations of adversarial learning and re-
188 gression losses [20, 26, 53], and combinations of flows
189 and VAEs [86]. Following the impressive performance of
190 text-prompted diffusion models for generating images [77]
191 and human motion [38, 87, 109], diffusion models have
192 seen rapid adoption for 3D gesture-motion generation. As
193 diffusion models require many neural-network evaluations
194 during synthesis, which is slow, flow matching [51] has
195 subsequently been investigated for faster synthesis of high
196 quality output, both for human motion [31, 62] and TTS
197 [25, 48, 63]. Similar to LLMs and large TTS models, recent
198 efforts have also wholly or partly modelled gestures autore-
199 gressively as a sequence of discrete tokens [64, 99, 107].

200 The most recent large-scale comparison of gesture-
201 generation models, the GENE Challenge 2023 [44], found
202 that the two strongest methods [17, 100] (which are exten-
203 sions of [7, 98]) were based on diffusion models. Among
204 these, [17] made use of self-supervised text-and speech em-
205 beddings from data2vec [8], subsequently aligned with ges-
206 ture motion using CLIP [72] training, to improve the co-
207 herence between gestures and the two speech-input modal-
208 ities. In addition to modelling beat gestures, the approach
209 recognises the need for additional input modalities to gen-
210 erate representational gestures, such as iconic and deictic
211 pointing [18], for more nuanced and contextually relevant
212 non-verbal communication.

213 Our data-synthesis pipeline leverages their approach to
214 create synthetic training gestures that well match the syn-
215 thetic speech text and audio input.

216 2.4. Joint synthesis of speech and gestures

217 Speech synthesis and gesture generation have traditionally
218 been treated as separate problems, performed on different
219 data by distinct research communities. TTS is mainly devel-

oped for read-aloud speech, whereas co-speech gesturing is
220 more closely associated with conversational settings. 221

222 Joint synthesis of speech and motion was first consid-
223 ered by [78]. The first neural model was DurIAN [106],
224 which simultaneously generated speech audio and 3D fac-
225 ial expressions, albeit for speech read aloud. [6] trained
226 separate deep-learning TTS and speech-to-gesture systems
227 to synthesise speech and 3D motion for the same speaker
228 and the same (spontaneous) speaking style. This was fol-
229 lowed by [94], which investigated adapting and extending
230 AR [79] and NAR [36] neural TTS models to perform joint
231 multimodal synthesis. Their joint models reduced the num-
232 ber of parameters needed over [6], but the best model (the
233 one based on [79]) required complex multi-stage training to
234 speak intelligibly and did not improve quality.

235 Diff-TTSG [61] advanced joint speech-and-gesture syn-
236 thesis by employing probabilistic modelling, specifically a
237 strong denoising probabilistic model (DPMs) [82] building
238 on the TTS work in [69]. This model could be trained on
239 speech-and-gesture data from scratch in one go and pro-
240 duced improved results over [94], but internally used sepa-
241 rate pipelines for producing the two output modalities, lead-
242 ing to suboptimal coherence between them. Match-TTSG
243 [62] improved on this aspect by using a compact and uni-
244 fied decoder to jointly sample both output modalities. It
245 also used conditional flow matching [51] rather than diffu-
246 sion, for much faster output synthesis. Experiments found
247 that Match-TTSG improved on the previous best model in
248 all respects, establishing it as the current state of the art.

249 Most of the above models were trained only on small,
250 parallel multimodal datasets from a single speaker. (The
251 one exception is [94], which required pre-training part of
252 the network on a TTS corpus to produce intelligible out-
253 put at all.) The results in [62] show that, e.g., the synthetic
254 speech falls short of human-level naturalness, and the qual-
255 ity we find from systems trained on very large datasets. Ac-
256 cordingly, we propose to circumvent the data limitation by
257 using strong unimodal synthesisers to create a large syn-
258 thetic training corpus for our joint model.

259 2.5. Training on synthetic data

260 The idea of training deep neural models on the output of
261 other such models has an extensive history. This was orig-
262 inally proposed for classifiers [29], but has subsequently
263 been adapted to generative models, e.g., for TTS [89]. Syn-
264 thesis (and synthetic data) is also appealing in scenarios
265 where real data is scarce or difficult to obtain, as demon-
266 strated in applications to human poses and motion [90, 108].
267 It also allows for the creation of diverse and controlled
268 datasets that can enable more accurate and versatile mod-
269 els [35]. We here propose to generalise such approaches by
270 chaining together multiple unimodal synthesisers, to enable
271 training multimodal speech-and-gesture models.

272 There may be a risk that the individual unimodal synthes- 321
273 isers in the proposed approach could fail to capture mutual 322
274 information that connects the modalities, since the differ- 323
275 ent synthesisers are likely to be trained on non-overlapping 324
276 data. This could in turn lead to synthesis artefacts and 325
277 failure to recreate correlations and dependencies between 326
278 modalities in systems trained on the final synthetic mul- 327
279 timodal corpus. However, recent theoretical and practical 328
280 results demonstrate that little [55] or no [52, 65] parallel 329
281 data may suffice for learning joint distributions of multi-
282 ple random variables (modalities). This suggests that train-
283 ing on corpora generated by synthesisers built from non-
284 overlapping material might not be as risky as it might seem.

285 3. Method

286 In this section we first describe our method for creating
287 wholly synthetic multimodal datasets for pre-training syn-
288 thesis models, followed by a description of our modifica-
289 tions to the Match-TTSG architecture to improve durations,
290 prosody control, and multi-speaker data.

291 3.1. Creating synthetic training data

292 Our pipeline for creating synthetic training data had the fol-
293 lowing main steps:

- 294 1. Generating written sentences in the style of conversa-
295 tional speech transcriptions.
- 296 2. Synthesising diverse speech audio from the text.
- 297 3. Validating/filtering the synthetic speech audio using au-
298 tomatic speech recognition, and aligning the input text
299 with the synthesised audio.
- 300 4. Synthesising gestures from the generated speech audio
301 files and their corresponding time-aligned text.

302 We provide more detail in the following subsections.

303 3.1.1 Text generation

304 The first step was to create text sentences that can form the
305 basis of synthesising multimodal data in a conversational
306 style. For this we utilised GPT-4 [2] and deliberate prompt-
307 ing. Specifically, we prompted the model with a list of 50
308 text transcriptions sentences from the training split [61] of
309 the Trinity Speech-Gesture Dataset II (TSGD2) [19, 21],
310 each enclosed in triple quotes, followed by a prompt re-
311 questing the model to produce 50 additional phrases in the
312 same style (including hesitations and disfluencies as seen in
313 the transcriptions) but ignoring the content. Further prompt-
314 ing then followed, to make the model generate additional
315 output based around different emotions and scenarios, so as
316 to obtain a more diverse material. The emotional categories
317 we provided were: disgust, sadness, fear, frustration, sur-
318 prise, excitement, happiness, confusion, and denial. Our
319 prompting often gave similar instructions multiple times,
320 since we found that such redundancy led to more realistic

output. The main instruction prompt and a number of ex- 321
ample continuations can be found in Appendix A. 322

323 We utilised the above procedure to generate a total of 600
324 phrases, each approximately 250 characters in length. We
325 found that limiting the length of the prompt helps prevent
326 issues with the subsequent speech synthesis, which shows
327 a tendency to produce unintelligible or confabulated output
328 when processing overly long utterances. The 600 generated
329 phrases will be shared in future revisions of the paper.

330 3.1.2 Speech generation

331 The next step was to synthesise speech audio from the 600
332 LLM-generated phrases. For this, we considered multi-
333 ple TTS systems capable of multi-speaker and spontaneous
334 speech synthesis, including Bark¹, XTTS [15], and Eleven-
335 Labs². However, Bark exhibited frequent confabulations
336 and unexpected changes in speaker identity within a single
337 utterance, which seemed problematic for learning to
338 maintain a consistent vocal identity. Although ElevenLabs
339 demonstrated high-quality output, its status as a non-open
340 source and proprietary solution led us to exclude it. Ul-
341 timately, we selected XTTS for generating our synthetic
342 speech dataset, due to it combining more consistent syn-
343 thesis with a research-permissible license. We limited each
344 synthesised utterance to at most 400 XTTS speech tokens,
345 since anything longer than that is virtually certain too long
346 for our prompts, and thus must contain confabulation or
347 gibberish speech. For everything else, default XTTS syn-
348 thesis hyperparameters were used. In the end, each syn-
349 thesised audio utterance was around 20–23 seconds long,
350 taking about half that time to synthesise.

351 In order to obtain more diverse data containing multiple
352 speakers, each of the 600 phrases was synthesised 16 times,
353 once in each of 16 different voices. These voices were se-
354 lected as a gender-balanced set (8 male and 8 female speak-
355 ers) from the VCTK corpus [97], and elicited from XTTS
356 by seeding the synthesis of each individual utterance with
357 the audio of longest VCTK utterance spoken by the rele-
358 vant speaker as an acoustic prompt. These prompting utter-
359 ances tended to be around 9 seconds long. In total, we thus
360 synthesised $16 \times 600 = 9600$ audio utterances.

361 Interestingly, despite the spontaneous nature of the in-
362 put phrases, we found that false starts and fillers explicitly
363 present in the input were sometimes omitted in the XTTS
364 output. This could be partly due to the choice of tempera-
365 ture parameter at synthesis time (the default, 0.65), which
366 favours more consistent and likely output, and partly due
367 to the public English-language training datasets cover read
368 rather than spontaneous speech. Since XTTS furthermore
369 was prompted using a snippet of read-aloud speech audio

¹<https://github.com/suno-ai/bark>

²<https://elevenlabs.io/>

370 from VCTK, the output audio tended to sound more like
371 reading than speaking spontaneously.

372 3.1.3 Data filtering and forced alignment

373 Following speech synthesis, a number of data-processing
374 steps were performed to obtain a suitable dataset for train-
375 ing a strong gesture-generation system. To begin with, all
376 synthesised audio utterances longer than 25 seconds were
377 immediately and permanently discarded, since these over-
378 whelmingly tended to contain issues related to confabula-
379 tion and the like. The output from XTTS did not have ex-
380 act fidelity to the text it was prompted with, so automatic
381 speech recognition (ASR) was used to get more accurate in-
382 put to the gesture-generation system. ASR was performed
383 using Whisper [73], using the `medium.en` model, which
384 has in previous uses proven to be less prone to confabula-
385 tion than the large variants, whilst providing sufficient accu-
386 racy. Interestingly, Whisper tended to prefer British English
387 spelling, possibly since VCTK was recorded in the UK. The
388 ASR derived transcripts then replaced the original TTS in-
389 put text for each utterance in all subsequent processing.

390 The gesture-generation system we chose for the final
391 synthesis ([17]) requires word-level timestamps for the text
392 transcriptions. Although we considered several tools that
393 attempt to obtain word timings from Whisper directly, none
394 were sufficiently accurate for our needs. Instead, we ob-
395 tained the requisite timings using the Montreal Forced
396 Aligner (MFA) [56]. Text input to MFA was processed
397 word-by-word to remove leading and trailing punctuation
398 and to perform case folding to lower case. Utterances that
399 MFA failed to align were also excluded from consideration.

400 Following the filtering and alignment process, we were
401 left with 8173 audio utterances for our final synthetic
402 dataset, meaning that 1427 utterances (about 15%) were
403 discarded during the filtering step. The remaining data had
404 a total duration of 37.6 hours, which also ended up being
405 the size of the final synthetic training corpus.

406 3.1.4 Gesture generation

407 We used a recent diffusion-based gesture-generation
408 method [17] that performed well in a large comparative
409 evaluation [44] to generate synthetic gesture data. That sys-
410 tem leveraged `data2vec` [8] embeddings to represent audio
411 input, which help achieve a more speaker-independent rep-
412 resentation. On top of that, [44] introduced a Contrastive
413 Speech and Motion Pretraining (CSMP) module, to learn
414 joint embeddings of speech and gesture that can strengthen
415 the semantic coupling between these modalities. By utilis-
416 ing the output of the CSMP module as a conditioning sig-
417 nal within the diffusion-based gesture-synthesis model, the
418 system can generate co-speech gestures that are human-like
419 and semantically aware, thereby improving the quality and

appropriateness of the generated gestures to the spoken con- 420
tent. The CSMP module requires word-level timestamps, 421
which is why forced-alignment was performed in Sec. 3.1.3. 422

423 Since this paper is focused on multimodal synthesis from
424 data where no interlocutor is present or recorded (i.e., not
425 back-and-forth conversations), interlocutor-related inputs
426 were removed from the architecture. The input is thus an
427 audio track with time-aligned text transcripts. We used the
428 pre-trained weights from [17] for the CSMP module and re-
429 trained the diffusion-based gesture model to comply with
430 the change of input, using the same architecture and learn-
431 ing rate as in the paper. The training was done using two
432 NVIDIA RTX3090 GPUs (194k updates, each with batch
433 size 60) on the subset of the Talking With Hands (TWH)
434 dataset [49] provided in the GENE 2023 Challenge [44].
435 We used the trained system to generate text-and-audio-
436 driven gestures for the 8173 previously transcribed syn-
437 thetic speech utterances, and used Autodesk MotionBuilder
438 after synthesis to retarget the output motion to the skele-
439 ton of the TSGD2 data and visualiser in Sec. 4.1. While
440 the synthesised motion encompasses the full body (without
441 fingers), we only consider upper-body motion in this work.
442 Compared to conventional conditioning approaches where
443 audio is represented using mel-spectrograms, the speaker-
444 independent `data2vec` embeddings in the CSMP module are
445 expected to better handle the differences between natural
446 and synthetic voices during synthesis, thus making it fea-
447 sible to generate large amounts of gesture data based on
448 synthetic speech without undue degradations due to domain
449 mismatch. This data was used to train the different multi-
450 modal synthesis systems considered in our experiments.

451 3.2. Proposed multimodal synthesis system

452 The current state of the art in joint speech-and-gesture syn-
453 thesis is Match-TTSG [62], a non-autoregressive model
454 which uses conditional flow matching (OT-CFM) [51] to
455 learn Ordinary Differential Equations (ODEs) with more
456 linear vector fields than continuous-time diffusion models
457 [82] create. Such simpler vector fields offer advantages for
458 easier learning and faster synthesis.

459 We extend the Match-TTSG framework in three ways:

- 460 1. Probabilistic instead of deterministic duration mod-
461 elling, which can benefit deep generative NAR TTS [37].
- 462 2. Additional prosody-prediction modules, which are
463 widely used in NAR TTS [75, 112].
- 464 3. A speaker-identity input, as necessary for pre-training on
465 the multispeaker data in the large synthetic training set.

466 We call the resulting system *MAGI* for *Multimodal Audio*
467 *and Gesture, Integrated*; see Fig. 2 for a diagram.

468 For (1), we augment the original Match-TTSG architec-
469 ture with a probabilistic duration predictor based on OT-
470 CFM, as introduced in [48], to learn distributions over
471 speech and gesture durations. This is trained jointly with

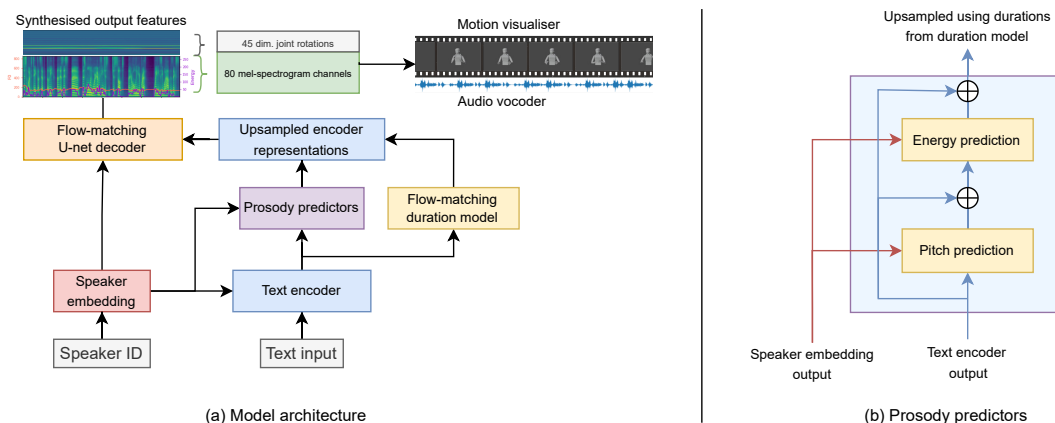


Figure 2. Schematic overview of the proposed MAGI architecture and its prosody predictor.

472 the rest of the system. It replaces the deterministic duration
473 predictor in Match-TTSG, inherited from [25, 36, 63, 69,
474 75, 112], and uses the same network architecture.

475 To learn better prosody correlations and enable control
476 over the output, we drew inspiration from [75, 112] and
477 incorporated two prosody-predictor modules into our system:
478 one for pitch prediction and one for energy prediction, both
479 using the same architecture and hyperparameters as the
480 *variance adaptor* in [75]. Such prosody predictors improve
481 the synthesis as they enable the model to learn a less over-
482 smoothed representation, thereby enhancing the variability
483 of the generated output by conditioning the synthesis process
484 on additional prosodic features [76]. The pitch of the training
485 data utterances was extracted using the PyWorld wrapper
486 for the WORLD vocoder³ with linear interpolation applied in
487 unvoiced segments to achieve continuous pitch contours for the
488 entire utterances. We employed a bucketing approach similar
489 to [75], separately for pitch and energy, to turn predicted
490 continuous values into embedding vectors to be summed with
491 the text-encoder output vectors. However, in contrast to [75],
492 we performed token-level prediction instead of frame-level
493 prediction for the two prosodic properties, since it has been
494 stated⁴ that this improves the synthesis whilst reducing
495 memory consumption.

496 Like in [69], Match-TTSG includes a projection layer
497 that maps the text-encoder output vectors onto a predicted
498 average output vector per token (sub-phone). These averages
499 are used for the so-called *prior loss* in the monotonic
500 alignment search. The process of sampling the output features
501 (i.e., the flow-matching decoder) is also conditioned on
502 these predicted average vectors. However, the latter can
503 introduce an information bottleneck, since averages do not
504 include information about variance, correlations, or higher
505 moments of the output distribution. To improve information
506 flow we instead condition the MAGI decoder directly on the

last layer of the text-encoder, prior to the projection layer. 507

508 Finally, we added a speaker embedding for multispeaker
509 synthesis. Specifically, we used a one-hot speaker vector
510 to represent the 16 different speakers in the synthetic training
511 data. This vector was concatenated to other inputs at
512 multiple stages of the synthesis process, including the text
513 encoder, prosody predictors and decoder. The idea with this
514 was to minimise information loss and ensure coherent output
515 across different speaker identities. Since the concatenated
516 vectors only have 16 elements, the impact on model
517 parameter count is small (an increase of a few thousand).

4. Experiments 518

519 This section experimentally compares our proposed training
520 method and architecture with the previous state-of-the-art
521 method Match-TTSG [62]. Since this is a synthesis work,
522 the gold standard approach to evaluation – and thus the focus
523 of our experimental validation – is subjective user studies.
524 The experiments closely follow those in previous joint
525 synthesis works [61, 62], which in turn follows established
526 practices in speech [32] and gesture evaluation [44].

4.1. Data and systems 527

528 To test the effectiveness of our method we carried out 3 different
529 subjective evaluations with systems trained on Trinity
530 Speech-Gesture Dataset II (TSGD2) [22], a dataset containing
531 6 hours of multimodal data: recordings of time-aligned
532 44.1 kHz audio coupled with 120 FPS marker-based 3D motion
533 capture, in which a male native speaker of Hiberno-
534 English discusses a variety of topics whilst gesturing freely.
535 The same train-test split of the data was used as in [61],
536 with around 4.5 hours of training data – much less than the
537 38 hours of synthetic multimodal data we created.

538 We trained Match-TTSG (**MAT**) containing 30.2M parameters,
539 and MAGI (**MAGI**) containing 31.6M parameters for 300k
540 steps on only the TSGD2 data, we refer to these conditions
541 **MAT-T** and **MAGI-T** respectively. We also took

³<https://pypi.org/project/pyworld/>

⁴<https://github.com/ming024/FastSpeech2?tab=readme-ov-file#implementation-issues>

the same two architectures (albeit with one-hot speaker vectors for Match-TTSG) and first pre-trained them for 200k updates on the synthetic multispeaker data, followed by fine-tuning for 100k updates on TSGD2. We refer to these as **MAT-FT** and **MAGI-FT**. Output samples for held-out sentences were synthesised using 100 neural function evaluations (NFEs; equivalent to number of Euler-forward steps used by the ODE solver) for audio-and-motion synthesis, whilst 10 NFEs were used for the preceding stochastic duration modelling, since it is lower-dimensional and converged more rapidly. Training and synthesis were performed on NVIDIA RTX 3090 GPUs with a batch size of 32.

15 utterances from the held-out set were used to evaluate each modality individually. We used pretrained Universal HiFi-GAN [39] to generate vocoded but otherwise natural speech referred to as **NAT**. We used the same vocoder to generate waveforms from the output mel spectrograms synthesised by the trained multimodal-synthesis systems, while Blender was used to render the motion representations into 3D avatar video, using exactly the same upper-body avatar and visualiser as in [61, 63]. The motion data was represented as rotational representation using exponential maps [24] of 45-dim pose vectors and were downsampled to 86.13 FPS using cubic interpolation to match the frame rate of the mel-spectrograms.

4.2. Evaluation setup

To gain an objective insight into the intelligibility of the synthetic speed, we synthesised the test set sentences from TSGD2, which we then passed to Whisper ASR, to use the Word Error Rate (WER) results as an indicator of their intelligibility. For subjective evaluation, user studies are the gold standard when evaluating synthesis methods. Following [61], we used comprehensive evaluation, conducting individual studies of each generated modality. We additionally evaluate the appropriateness of the modalities in terms of each other, to determine how well they fit together.

In our studies, participants had an interface with five unique response choices, with the exact details varying slightly across different investigations. All participants were native English speakers recruited through the Prolific⁵ crowdsourcing platform. Each test was designed to last around 20 minutes and participants were compensated 4 GBP (12 GBP/hr) for participation. For the purpose of statistical examination, we converted responses into numerical values. These values were then analysed for statistical significance at the 0.05 threshold using pairwise t-tests.

4.2.1 Speech-quality evaluation

To assess perceived naturalness of the synthesized speech, we employed the Mean Opinion Score (MOS) testing ap-

proach, drawing inspiration from the Blizzard Challenge for text-to-speech systems [70]. Participants were asked, “How natural does the synthesized speech sound?”, rating their responses on a scale from 1 to 5, where 1 represented “Completely unnatural” and 5 indicated “Completely natural.” The intermediary values of 2 to 4 were provided without textual descriptions. Each participant evaluated 15 stimuli per system and 4 attention checks resulting in a total of 525 responses per condition by 35 participants. Fine-tuning with synthetic data led to performance enhancements for both MAGI and MAT, reducing the WER from 13.28% in MAGI-T to 9.29% in MAGI-FT, and from 12.26% in MAT-T to 8.35% in MAT-FT.

4.2.2 Motion-quality evaluation

We evaluate motion quality using video stimuli that only visualised motion, without any audio, in order to have an independent assessment of motion quality. This ensures that ratings are not affected by speech and follows the practice of recent evaluations of gesture quality [33, 74]. Similarly to the speech evaluation, participants were asked “How natural and humanlike the gesture motion appear?”, and gave responses on a scale of 1 (“Completely unnatural”) to 5 (“Completely natural”). The number of stimuli and attention checks were identical to the speech-only evaluation.

4.2.3 Speech-and-motion appropriateness evaluation

We finally evaluated how appropriate the generated speech and motion were for each other, whilst controlling for the effect of their individual quality following [33, 45, 62, 74, 105]. For each speech segment and condition, we created two video stimuli: one with the original video and sound, and the other combining the original speech audio with motion from a different video clip, adjusting the motion speed to align with the audio duration. Both videos feature comparable motion quality and characteristics from the same condition, but only one video’s motion is synchronised with the audio track, without indicating which video is which.

The test inquired which character’s motion most accurately matched the speech in rhythm, intonation, and meaning. Participant ability to identify the correctly synchronised video indicates a strong rhythmic and/or semantic link between generated motion and speech. Following [61] we opted for five response choices instead of the typical three for better resolution. Options were “Left is much better”, “Left is slightly better”, “Both are equal”, “Right is slightly better”, “Right is much better”. For the purposes of analysis, codes in the range of -2 to 2 were assigned to each response, as in [61], with -2 representing the participant’s preference for the mismatched stimulus and 2 the matched stimulus. Participants reviewed motions from 14 of the 15 segments, displayed as 7 screens of pairs of videos, plus

⁵<https://www.prolific.com/>

Table 1. Result of three evaluations showing Mean Opinion Scores (MOS) and 95% confidence intervals.

Condition	Speech	Gesture	Speech & Gesture
NAT	4.30±0.06	4.10±0.08	1.10±0.10
MAT-T	3.43±0.10	3.28±0.11	0.52±0.10
MAT-FT	3.56±0.10	3.39±0.09	0.56±0.09
MAGI-T	3.44±0.09	3.11±0.10	0.51±0.09
MAGI-FT	3.62±0.08	3.52±0.11	0.60±0.09

two audio and two video attention checks, covering all conditions for these segments. 70 people completed the test, yielding 490 responses per system.

5. Results and discussion

Our investigation revealed several key insights into the effect of pre-training and architectural modifications. Pre-training on synthetic data markedly enhanced the quality of synthesised speech, though adjustments to the architecture did not significantly alter its naturalness. Despite this, both MAGI-FT and MAT-FT yielded higher Mean Opinion Scores (MOS), albeit without statistical significance. Notably, the MAGI facilitated greater control over pitch and energy—a feature absent in the original MAT framework. However, despite improvements, the synthesised speech did not achieve the level of naturalness present in the human-recorded speech from the held-out set, see Table 1.

In terms of synthesised gestures, MAGI outperformed other conditions in human-likeness. However, they remained inferior to human-motion reference data. The influence of synthetic data pre-training and the proposed model’s architecture on gesture synthesis presented a more nuanced picture. Specifically, pre-training on synthetic data only significantly benefited the proposed model, and, intriguingly, the MAGI enhanced gestures in a larger dataset but had the opposite effect on a smaller dataset. This discrepancy might stem from the prosody predictors in our model being trained on per-phone rather than per-frame data, leading to a scarcity of training data for these predictors in smaller datasets. However, with adequate pre-training on expansive datasets, these models demonstrated better convergence. These findings align with prior speech evaluations, where the novel architecture’s advantages were more pronounced following pre-training on a larger dataset.

Further, no model matched the cross-modal appropriateness found in multimodal human recordings, echoing the challenges observed in unimodal gesture synthesis where recent evaluations did not approach the appropriateness of human data [45, 105]. Although MAGI, pre-trained on synthetic data, showcased superior performance, it did not significantly exceed the existing benchmarks in synthesis systems. This observation may be attributed to the inher-

ent difficulty in discerning significant differences in appropriateness, as opposed to naturalness or human-likeness, and the comparison against a robust baseline without alterations that directly influence cross-modal synthesis aspects. Lastly, the accuracy of capturing cross-modal aspects might be least represented in synthetic datasets created from unimodal synthesizers trained on non-cohesive data.

5.1. Pitch and energy control

As stated, the proposed multi-stage architecture with separate prosody predictors allows for modifying or substituting the pitch and energy contours before synthesis. This enables direct control of prosodic properties of the speech, with the synthesis process having the option to adjust the gestures to match. On our anonymous webpage cvprhuman24.github.io/MAGI we provide example videos showing the effect that modifying (scaling) the pitch and energy contours returned by the predictors has on the synthesised output. One can observe that reducing the pitch seems to promote creaky voice, which makes sense from a speech-production perspective and fits earlier findings from autoregressive TTS on spontaneous-speech data [47].

6. Conclusion and future work

We have described improvements to the joint and simultaneous multimodal synthesis of speech audio and 3D gesture motion from text. Specifically, we propose pre-training on data synthesised by a chain of strong unimodal synthesis systems to address the shortage of multimodal training data. We also augment the state-of-the-art architecture for speech-and-gesture synthesis, Match-TTSG, with a stochastic duration model, TTS-inspired prosody predictors for controllability, and the ability to perform multi-speaker synthesis. The final model, called Multimodal Audio and Gesture, Integrated (MAGI), is radically smaller than those that generated the synthetic data. Experiments confirm that pre-training on synthetic data significantly improved unimodal speech and gesture quality. The architectural improvements reaped benefits when pre-training on large amounts of synthetic data, with the added prosody control having a clear effect on the audio output.

Relevant future work includes investigating alternative options for mitigating the shortage of multimodal training data, such as pre-training on data lacking one or more of the modalities, incorporating RL-based approaches, particularly effective for generation of situated gestures as in [18], or (following the CSMP methodology [17]) leveraging various self-supervised representations trained on large amounts of data. Possible architectural extensions including flow matching for pitch and energy, and similar control over motion properties such as gesture radius and symmetry [5].

732 **References**

- 733 [1] Yelaman Abdullin, Diego Molla-Aliod, Bahadorreza
734 Ofoghi, John Yearwood, and Qingyang Li. Synthetic di-
735 alogue dataset generation using llm agents. *arXiv preprint*
736 *arXiv:2401.17461*, 2024. 2
- 737 [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ah-
738 mad, Ilge Akkaya, Florencia Leoni Aleman, et al. GPT-4
739 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
740 1, 2, 4
- 741 [3] Chaitanya Ahuja, Dong Won Lee, Ryo Ishii, and Louis-
742 Philippe Morency. No gestures left behind: Learning rela-
743 tionships between spoken language and freeform gestures.
744 In *Proc. EMNLP*, pages 1884–1895, 2020. 1
- 745 [4] Simon Alexanderson. The StyleGestures entry to the GE-
746 NEA Challenge 2020. In *Proc. GENEA Workshop*, 2020.
747 3
- 748 [5] Simon Alexanderson, Gustav Eje Henter, Taras
749 Kucherenko, and Jonas Beskow. Style-controllable
750 speech-driven gesture synthesis using normalising flows.
751 *Comput. Graph. Forum*, 39(2):487–496, 2020. 3, 8
- 752 [6] Simon Alexanderson, Éva Székely, Gustav Eje Henter,
753 Taras Kucherenko, and Jonas Beskow. Generating coherent
754 spontaneous speech and gesture from text. In *Proc. IVA*,
755 pages 1–3, 2020. 1, 3
- 756 [7] Simon Alexanderson, Rajmund Nagy, Jonas Beskow, and
757 Gustav Eje Henter. Listen, denoise, action! Audio-driven
758 motion synthesis with diffusion models. *ACM Trans.*
759 *Graph.*, 42(4):1–20, 2023. 3
- 760 [8] Alexei Baevski, Wei-Ning Hsu, Qiantong Xu, Arun Babu,
761 Jiatao Gu, and Michael Auli. data2vec: A general frame-
762 work for self-supervised learning in speech, vision and lan-
763 guage. In *Proceedings of the International Conference on*
764 *Machine Learning*, pages 1298–1312, 2022. 3, 5
- 765 [9] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell,
766 Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort,
767 Deep Ganguli, Tom Henighan, et al. Training a helpful and
768 harmless assistant with reinforcement learning from human
769 feedback. *arXiv preprint arXiv:2204.05862*, 2022. 2
- 770 [10] James Betker. Better speech synthesis through scaling.
771 *arXiv preprint arXiv:2305.07243*, 2023. 2
- 772 [11] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng
773 Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee,
774 Yufei Guo, et al. Improving image generation with better
775 captions, 2023. 2
- 776 [12] Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene
777 Kharitonov, Olivier Pietquin, Matt Sharifi, Dominik Rob-
778 lek, Olivier Teboul, David Grangier, Marco Tagliasacchi,
779 et al. Audioldm: a language modeling approach to audio
780 generation. *IEEE/ACM Transactions on Audio, Speech, and*
781 *Language Processing*, 2023. 2
- 782 [13] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Sub-
783 biah, Jared D. Kaplan, Prafulla Dhariwal, et al. Language
784 models are few-shot learners. In *Advances in Neural In-*
785 *formation Processing Systems*, pages 1877–1901, 2020. 1,
786 2
- 787 [14] Justine Cassell, Joseph Sullivan, Scott Prevost, and Eliz-
abeth Churchill. *Embodied conversational agents*. MIT
press, 2000. 1
- [15] Coqui.ai. xttts - tts 0.22.0 documentation, 2023. 2, 4
- [16] Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and
Yossi Adi. High fidelity neural audio compression. *arXiv*
preprint arXiv:2210.13438, 2022. 2
- [17] Anna Deichler, Shivam Mehta, Simon Alexanderson, and
Jonas Beskow. Diffusion-based co-speech gesture genera-
tion using joint text and audio representation. In *Proceed-*
ings of the 25th International Conference on Multimodal
Interaction, pages 755–762, 2023. 3, 5, 8
- [18] Anna Deichler, Siyang Wang, Simon Alexanderson, and
Jonas Beskow. Learning to generate pointing gestures
in situated embodied conversational agents. *Frontiers in*
Robotics and AI, 10:1110534, 2023. 3, 8
- [19] Ylva Ferstl and Rachel McDonnell. Investigating the use of
recurrent motion modelling for speech gesture generation.
In *Proc. IVA*, pages 93–98, 2018. 4
- [20] Ylva Ferstl and Rachel McDonnell. Multi-task learning for
continuous control of non-verbal behaviour in humanoid
social robots. In *Proceedings of the ACM/IEEE Inter-*
national Conference on Human-Robot Interaction (HRI),
pages 411–420. IEEE, 2019. 3
- [21] Ylva Ferstl, Michael Neff, and Rachel McDonnell. Ad-
versarial gesture generation with realistic gesture phasing.
Comput. Graph., 89:117–130, 2020. 4
- [22] Ylva Ferstl, Michael Neff, and Rachel McDonnell. Ex-
pressGesture: Expressive gesture generation from speech
through database matching. *Comput. Animat. Virt. W.*, page
e2016, 2021. 2, 6
- [23] Saeed Ghorbani, Ylva Ferstl, Daniel Holden, Nikolaus F.
Troje, and Marc-André Carbonneau. ZeroEGGS: Zero-shot
example-based gesture generation from speech. *Comput.*
Graph. Forum, 42(1):206–216, 2023. 3
- [24] F. Sebastian Grassia. Practical parameterization of rotations
using the exponential map. *J. Graph. Tool.*, 3(3):29–48,
1998. 7
- [25] Yiwei Guo, Chenpeng Du, Ziyang Ma, Xie Chen, and Kai
Yu. VoiceFlow: Efficient text-to-speech with rectified flow
matching. In *Proc. ICASSP*, 2024. 2, 3, 6
- [26] Ikhsanul Habibie, Weipeng Xu, Dushyant Mehta, Lingjie
Liu, Hans-Peter Seidel, Gerard Pons-Moll, Mohamed El-
gharib, and Christian Theobalt. Learning speech-driven
3D conversational gestures from video. In *Proceedings of*
the International Conference on Intelligent Virtual Agents,
pages 101–108, 2021. 3
- [27] Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen,
Christopher Hesse, Jacob Jackson, Heewoo Jun, Tom B.
Brown, Prafulla Dhariwal, Scott Gray, et al. Scaling laws
for autoregressive generative modeling. *arXiv preprint*
arXiv:2010.14701, 2020. 1
- [28] Laura Birka Hensel, Nutchanon Yongsatianchot, Parisa
Torshizi, Elena Minucci, and Stacy Marsella. Large lan-
guage models in textual analysis for gesture selection. In
Proceedings of the 25th International Conference on Mul-
timodal Interaction, pages 378–387, 2023. 2

- [29] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. In *NIPS Deep Learning and Representation Learning Workshop*, 2015. 3
- [30] Autumn B Hostetter. When do gestures communicate? a meta-analysis. *Psychological Bulletin*, 133(2):297, 2007. 1
- [31] Vincent Tao Hu, Wenzhe Yin, Pingchuan Ma, Yunlu Chen, Basura Fernando, Yuki M Asano, Efstratios Gavves, Pascal Mettes, Björn Ommer, and Cees G. M. Snoek. Motion flow matching for human motion synthesis and editing. *arXiv preprint arXiv:2312.08895*, 2023. 3
- [32] ITU-T P.800. Methods for subjective determination of transmission quality. Standard, ITU, 1996. 6
- [33] Patrik Jonell, Taras Kucherenko, Gustav Eje Henter, and Jonas Beskow. Let’s face it: Probabilistic multi-modal interlocutor-aware generation of facial gestures in dyadic settings. In *Proc. IVA*, 2020. 7
- [34] Adam Kendon. How gestures can become like words. In *Cross-Cultural Perspectives in Nonverbal Communication*. C. J. Hogrefe, Inc., 1988. 1
- [35] Hyeonwoo Kim, Pablo Garrido, Ayush Tewari, Weipeng Xu, Justus Thies, Matthias Niessner, Patrick Pérez, Christian Richardt, Michael Zollhöfer, and Christian Theobalt. Neural style-preserving visual dubbing. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 2535–2545, 2019. 3
- [36] Jaehyeon Kim, Sungwon Kim, Jungil Kong, and Sungroh Yoon. Glow-TTS: A generative flow for text-to-speech via monotonic alignment search. In *Proc. NeurIPS*, pages 8067–8077, 2020. 2, 3, 6
- [37] Jaehyeon Kim, Jungil Kong, and Juhee Son. VITS: Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech. In *Proc. ICML*, pages 5530–5540, 2021. 2, 3, 5
- [38] Jihoon Kim, Jiseob Kim, and Sungjoon Choi. Flame: Free-form language-based motion synthesis & editing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 8255–8263, 2023. 3
- [39] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. HiFi-GAN: Generative adversarial networks for efficient and high fidelity speech synthesis. In *Proc. NeurIPS*, pages 17022–17033, 2020. 7
- [40] Jungil Kong, Jihoon Park, Beomjeong Kim, Jeongmin Kim, Dohee Kong, and Sangjin Kim. VITS2: Improving quality and efficiency of single-stage text-to-speech with adversarial learning and architecture design. In *Proc. Interspeech*, pages 4374–4378, 2023. 3
- [41] Stefan Kopp and Ipke Wachsmuth. Synthesizing multi-modal utterances for conversational agents. In *Computer Animation and Virtual Worlds*, pages 39–52. Wiley Online Library, 2004. 1
- [42] Taras Kucherenko, Patrik Jonell, Sanne van Waveren, Gustav Eje Henter, Simon Alexanderson, Iolanda Leite, and Hedvig Kjellström. Gesticulator: A framework for semantically-aware speech-driven gesture generation. In *Proceedings of the ACM International Conference on Intelligent Virtual Agents*, pages 242–250, 2020. 1
- [43] Taras Kucherenko, Rajmund Nagy, Michael Neff, Hedvig Kjellström, and Gustav Eje Henter. Multimodal analysis of the predictability of hand-gesture properties. In *Proceedings of the International Conference on Autonomous Agents and Multiagent Systems*, pages 770–779, 2022. 1
- [44] Taras Kucherenko, Rajmund Nagy, Youngwoo Yoon, Jieyeon Woo, Teodor Nikolov, Mihail Tsakov, and Gustav Eje Henter. The GENE Challenge 2023: A large-scale evaluation of gesture generation models in monadic and dyadic settings. In *Proceedings of the International Conference on Multimodal Interaction*, pages 792–801, 2023. 2, 3, 5, 6
- [45] Taras Kucherenko, Pieter Wolfert, Youngwoo Yoon, Carla Viegas, Teodor Nikolov, Mihail Tsakov, and Gustav Eje Henter. Evaluating gesture-generation in a large-scale open challenge: The GENE Challenge 2022. *arXiv preprint arXiv:2303.08737*, 2023. 7, 8
- [46] Mateusz Lajszczak, Guillermo Cambara Ruiz, Yang Li, Fatih Beyhan, Arent van Korlaar, Fan Yang, et al. Base tts: Lessons from building a billion-parameter text-to-speech model on 100k hours of data. *arXiv*, 2024. 2
- [47] Harm Lameris, Shivam Mehta, Gustav Eje Henter, Joakim Gustafson, and Éva Székely. Prosody-controllable spontaneous TTS with neural HMMs. In *Proc. ICASSP*, 2023. 3, 8
- [48] Matthew Le, Apoorv Vyas, Bowen Shi, Brian Karrer, Leda Sari, Rashel Moritz, Mary Williamson, Vimal Manohar, Yossi Adi, Jay Mahadeokar, et al. Voicebox: Text-guided multilingual universal speech generation at scale. *arXiv preprint arXiv:2306.15687*, 2023. 2, 3, 5
- [49] Gilwoo Lee, Zhiwei Deng, Shugao Ma, Takaaki Shiratori, Siddhartha S. Srinivasa, and Yaser Sheikh. Talking With Hands 16.2 M: A large-scale dataset of synchronized body-finger motion and audio for conversational motion analysis and synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 763–772, 2019. 2, 5
- [50] Naihuan Li, Shujie Liu, Yanqing Liu, Sheng Zhao, and Ming Liu. Neural speech synthesis with transformer network. In *Proceedings of the AAAI conference on artificial intelligence*, pages 6706–6713, 2019. 2
- [51] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. In *Proc. ICLR*, 2023. 1, 3, 5
- [52] Alexander H. Liu, Cheng-I Jeff Lai, Wei-Ning Hsu, Michael Auli, Alexei Baevski, and James Glass. Simple and effective unsupervised speech synthesis. In *Proc. Interspeech*, pages 843–847, 2022. 4
- [53] Haiyang Liu, Zihao Zhu, Naoya Iwamoto, Yichen Peng, Zhengqing Li, You Zhou, Elif Bozkurt, and Bo Zheng. BEAT: A large-scale semantic and emotional multi-modal dataset for conversational gestures synthesis. In *Proceedings of the European Conference on Computer Vision*, pages 612–630, 2022. 2, 3
- [54] Yu Liu, Gelareh Mohammadi, Yang Song, and Wafa Johal. Speech-based gesture generation for robots and embodied agents: A scoping review. In *Proceedings of the International Conference on Human-Agent Interaction*, pages 31–38, 2021. 1

- [55] Soroosh Mariooryad, Matt Shannon, Siyuan Ma, Tom Bagby, David Kao, Daisy Stanton, Eric Battenberg, and RJ Skerry-Ryan. Learning the joint distribution of two sequences using little or no paired data. *arXiv preprint arXiv:2212.03232*, 2022. 4
- [56] Michael McAuliffe, Michaela Socolof, Sarah Mihuc, Michael Wagner, and Morgan Sonderegger. Montreal Forced Aligner: Trainable text-speech alignment using Kaldi. In *Proc. Interspeech 2017*, pages 498–502, 2017. 5
- [57] David McNeill. *Hand and mind: What gestures reveal about thought*. University of Chicago Press, 1992. 1
- [58] David McNeill. *Gesture and Thought*. University of Chicago Press, 2008. 1
- [59] Shivam Mehta, Éva Székely, Jonas Beskow, and Gustav Eje Henter. Neural HMMs are all you need (for high-quality attention-free TTS). In *Proc. ICASSP*, pages 7457–7461, 2022. 2
- [60] Shivam Mehta, Ambika Kirkland, Harm Lameris, Jonas Beskow, Éva Székely, and Gustav Eje Henter. OverFlow: Putting flows on top of neural transducers for better TTS. In *Proc. Interspeech*, 2023. 2
- [61] Shivam Mehta, Siyang Wang, Simon Alexanderson, Jonas Beskow, Éva Székely, and Gustav Eje Henter. Diff-TTSG: Denoising probabilistic integrated speech and gesture synthesis. In *Proc. SSW*, 2023. 1, 3, 4, 6, 7
- [62] Shivam Mehta, Ruibo Tu, Simon Alexanderson, Jonas Beskow, Éva Székely, and Gustav Eje Henter. Unified speech and gesture synthesis using flow matching. In *Proc. ICASSP*, 2024. 1, 2, 3, 5, 6, 7
- [63] Shivam Mehta, Ruibo Tu, Jonas Beskow, Éva Székely, and Gustav Eje Henter. Matcha-TTS: A fast TTS architecture with conditional flow matching. In *Proc. ICASSP*, 2024. 2, 3, 6, 7
- [64] Evonne Ng, Javier Romero, Timur Bagautdinov, Shaojie Bai, Trevor Darrell, Angjoo Kanazawa, and Alexander Richard. From audio to photoreal embodiment: Synthesizing humans in conversations. *arXiv preprint arXiv:2401.01885*, 2024. 3
- [65] Junrui Ni, Liming Wang, Heting Gao, Kaizhi Qian, Yang Zhang, Shiyu Chang, and Mark Hasegawa-Johnson. Unsupervised text-to-speech synthesis by unsupervised automatic speech recognition. In *Proc. Interspeech*, pages 461–465, 2022. 4
- [66] Simbarashe Nyatsanga, Taras Kucherenko, Chaitanya Ahuja, Gustav Eje Henter, and Michael Neff. A comprehensive review of data-driven co-speech gesture generation. *Comput. Graph. Forum*, 2023. 1, 3
- [67] Sewade Ogun, Vincent Colotte, and Emmanuel Vincent. Stochastic pitch prediction improves the diversity and naturalness of speech in Glow-TTS. In *Proc. Interspeech*, 2023. 2
- [68] Catherine Pelachaud, Norman I Badler, and Mark Steedman. Modeling and animating conversational agents. In *Adaptive hypertext and hypermedia*, pages 21–30. Springer, 1996. 1
- [69] Vadim Popov, Ivan Vovk, Vladimir Gogoryan, Tasnima Sadekova, and Mikhail Kudinov. Grad-TTS: A diffusion probabilistic model for text-to-speech. In *Proc. ICML*, pages 8599–8608, 2021. 2, 3, 6
- [70] Kishore Prahallad, Anandaswarup Vadapalli, Naresh Elluru, Gautam Mantena, Bhargav Pulugundla, Peri Bhaskararao, Hema A. Murthy, Simon King, Vasilis Karaiskos, and Alan W. Black. The Blizzard Challenge 2013–Indian language task. In *Proceedings of the Blizzard Challenge Workshop*, 2013. 7
- [71] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training, 2018. 2
- [72] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *Proceedings of the International Conference on Machine Learning*, pages 8748–8763, 2021. 3
- [73] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *Proceedings of the International Conference on Machine Learning*, pages 28492–28518, 2023. 5
- [74] Manuel Rebol, Christian Güti, and Krzysztof Pietroszek. Passing a non-verbal Turing test: Evaluating gesture animations generated from speech. In *Proceedings of the IEEE Conference on Virtual Reality and 3D User Interfaces*, pages 573–581, 2021. 7
- [75] Yi Ren, Chenxu Hu, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. FastSpeech 2: Fast and high-quality end-to-end text to speech. In *Proc. ICLR*, 2021. 2, 5, 6
- [76] Yi Ren, Xu Tan, Tao Qin, Zhou Zhao, and Tie-Yan Liu. Revisiting over-smoothness in text to speech. In *Proc. ACL*, pages 8197–8213, 2022. 6
- [77] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proc. CVPR*, pages 10684–10695, 2022. 1, 3
- [78] Maha Salem, Stefan Kopp, Ipke Wachsmuth, and Frank Joublin. Towards an integrated model of speech and gesture production for multi-modal robot behavior. In *Proc. RO-MAN*, pages 614–619, 2010. 1, 3
- [79] Jonathan Shen, Ruoming Pang, Ron J. Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng Chen, et al. Natural TTS synthesis by conditioning WaveNet on mel spectrogram predictions. In *Proc. ICASSP*, pages 4779–4783, 2018. 2, 3
- [80] Kai Shen, Zeqian Ju, Xu Tan, Yanqing Liu, Yichong Leng, Lei He, Tao Qin, et al. NaturalSpeech 2: Latent diffusion models are natural and zero-shot speech and singing synthesizers. *arXiv preprint arXiv:2304.09116*, 2023. 1, 2
- [81] Hubert Siuzdak. Vocos: Closing the gap between time-domain and fourier-based neural vocoders for high-quality audio synthesis. *arXiv preprint arXiv:2306.00814*, 2023. 2
- [82] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *Proc. ICLR*, 2021. 3, 5

- 1074 [83] Xu Tan, Tao Qin, Frank Soong, and Tie-Yan Liu. A
1075 survey on neural speech synthesis. *arXiv preprint*
1076 *arXiv:2106.15561*, 2021. 1, 2
- 1077 [84] Xu Tan, Jiawei Chen, Haohe Liu, Jian Cong, Chen Zhang,
1078 Yanqing Liu, Xi Wang, et al. NaturalSpeech: End-to-end
1079 text to speech synthesis with human-level quality. *arXiv*
1080 *preprint arXiv:2205.04421*, 2022. 1, 2
- 1081 [85] Paul Taylor. *Text-to-speech synthesis*. Cambridge Univer-
1082 sity Press, 2009. 1
- 1083 [86] Sarah Taylor, Jonathan Windle, David Greenwood, and Iain
1084 Matthews. Speech-driven conversational agents using con-
1085 ditional Flow-VAEs. In *Proceedings of the ACM Euro-
1086 pean Conference on Visual Media Production*, pages 6:1–
1087 6:9, 2021. 3
- 1088 [87] Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir,
1089 Daniel Cohen-Or, and Amit H. Bermano. Human motion
1090 diffusion model. In *Proceedings of the International Con-
1091 ference on Learning Representations*, 2023. 3
- 1092 [88] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier
1093 Martinet, Marie-Anne Lachaux, Timothée Lacroix, Bap-
1094 tiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar,
1095 et al. Llama: Open and efficient foundation language mod-
1096 els. *arXiv preprint arXiv:2302.13971*, 2023. 2
- 1097 [89] Aaron van den Oord, Yazhe Li, Igor Babuschkin, Karen
1098 Simonyan, Oriol Vinyals, Koray Kavukcuoglu, George
1099 Driessche, Edward Lockhart, Luis Cobo, Florian Stimberg,
1100 et al. Parallel WaveNet: Fast high-fidelity speech synthesis.
1101 In *Proceedings of the International Conference on Machine*
1102 *Learning*, pages 3918–3926, 2018. 3
- 1103 [90] Gül Varol, Javier Romero, Xavier Martin, Naureen Mah-
1104 mood, Michael J Black, Ivan Laptev, and Cordelia Schmid.
1105 Learning from synthetic humans. In *Proceedings of the*
1106 *IEEE Conference on Computer Vision and Pattern Recog-
1107 nition*, pages 109–117, 2017. 3
- 1108 [91] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
1109 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser,
1110 and Illia Polosukhin. Attention is all you need. *Advances*
1111 *in neural information processing systems*, 30, 2017. 1, 2
- 1112 [92] Petra Wagner, Zofia Malisz, and Stefan Kopp. Gesture and
1113 speech in interaction: An overview. *Speech Communica-
1114 tion*, 57:209–232, 2014. 1
- 1115 [93] Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang,
1116 Long Zhou, Shujie Liu, Zhuo Chen, Yanqing Liu, Huam-
1117 ing Wang, Jinyu Li, et al. Neural codec language models
1118 are zero-shot text to speech synthesizers. *arXiv preprint*
1119 *arXiv:2301.02111*, 2023. 2
- 1120 [94] Siyang Wang, Simon Alexanderson, Joakim Gustafson,
1121 Jonas Beskow, Gustav Eje Henter, and Éva Székely. Inte-
1122 grated speech and gesture synthesis. In *Proc. ICMI*, pages
1123 177–185, 2021. 1, 3
- 1124 [95] Bowen Wu, Chaoran Liu, Carlos T. Ishi, and Hiroshi Ishig-
1125 uro. Modeling the conditional distribution of co-speech
1126 upper body gesture jointly using conditional-GAN and
1127 unrolled-GAN. *Electronics*, 10(3):228, 2021. 3
- 1128 [96] Bowen Wu, Chaoran Liu, Carlos T. Ishi, and Hiroshi Ishig-
1129 uro. Probabilistic human-like gesture synthesis from speech
1130 using GRU-based WGAN. In *Companion Publication of*
the International Conference on Multimodal Interaction,
pages 194–201, 2021. 3
- [97] Junichi Yamagishi, Christophe Veaux, and Kirsten Mac-
Donald. CSTR VCTK corpus: English multi-speaker cor-
pus for CSTR voice cloning toolkit (version 0.92), 2019.
4
- [98] Sicheng Yang, Zhiyong Wu, Minglei Li, Zhensong Zhang,
Lei Hao, Weihong Bao, Ming Cheng, and Long Xiao. Dif-
fusestylegesture: stylized audio-driven co-speech gesture
generation with diffusion models. In *Proceedings of the*
International Joint Conference on Artificial Intelligence,
pages 5860–5868, 2023. 3
- [99] Sicheng Yang, Zhiyong Wu, Minglei Li, Zhensong Zhang,
Lei Hao, Weihong Bao, and Haolin Zhuang. QPGesture:
Quantization-based and phase-guided motion matching for
natural speech-driven gesture generation. In *Proceedings of*
the IEEE/CVF Conference on Computer Vision and Pattern
Recognition, pages 2321–2330, 2023. 3
- [100] Sicheng Yang, Haiwei Xue, Zhensong Zhang, Minglei Li,
Zhiyong Wu, Xiaofei Wu, Songcen Xu, and Zonghong Dai.
The diffusestylegesture+ entry to the genea challenge 2023.
In *Proceedings of the International Conference on Multi-
modal Interaction*, pages 779–785, 2023. 3
- [101] Yusuke Yasuda, Xin Wang, and Junichi Yamagishi. Ef-
fect of choice of probability distribution, randomness, and
search methods for alignment modeling in sequence-to-
sequence text-to-speech synthesis using hard alignment.
In *ICASSP 2020-2020 IEEE International Conference on*
Acoustics, Speech and Signal Processing (ICASSP), pages
6724–6728. IEEE, 2020. 2
- [102] Payam Jome Yazdian, Mo Chen, and Angelica Lim. Ges-
ture2Vec: Clustering gestures using representation learn-
ing methods for co-speech gesture generation. In *Proceed-
ings of the IEEE/RSJ International Conference on Intelli-
gent Robots and Systems*, 2022. 3
- [103] Hongwei Yi, Hualin Liang, Yifei Liu, Qiong Cao, Yandong
Wen, Timo Bolkart, Dacheng Tao, and Michael J Black.
Generating holistic 3D human motion from speech. In *Pro-
ceedings of the IEEE/CVF Conference on Computer Vision*
and Pattern Recognition, pages 469–480, 2023. 3
- [104] Youngwoo Yoon, Bok Cha, Joo-Haeng Lee, Minsu Jang,
Jaeyeon Lee, Jaehong Kim, and Geehyuk Lee. Speech ges-
ture generation from the trimodal context of text, audio, and
speaker identity. *ACM T. Graphic.*, 39(6):222:1–222:16,
2020. 1
- [105] Youngwoo Yoon, Pieter Wolfert, Taras Kucherenko, Carla
Viegas, Teodor Nikolov, Mihail Tsakov, and Gustav Eje
Henter. The GENE Challenge 2022: A large evaluation of
data-driven co-speech gesture generation. In *Proceedings*
of the International Conference on Multimodal Interaction,
pages 736–747, 2022. 7, 8
- [106] Chengzhu Yu, Heng Lu, Na Hu, Meng Yu, Chao Weng, Kun
Xu, Peng Liu, et al. DurIAN: Duration informed attention
network for multimodal synthesis. In *Proc. Interspeech*,
pages 2027–2031, 2020. 3
- [107] Xin Zeng, Xiaoyu Wang, Tengxiang Zhang, Chun Yu,
Shengdong Zhao, and Yiqiang Chen. GestureGPT: 1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

- 1188 Zero-shot interactive gesture understanding and ground-
1189 ing with large language model agents. *arXiv preprint*
1190 *arXiv:2310.12821*, 2023. 3
- [108] He Zhang, Sebastian Starke, Taku Komura, and Jun Saito.
1191 Motion synthesis and editing in low-dimensional spaces.
1192 *Computer Graphics Forum*, 39(8):509–521, 2020. 3
- [109] Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou
1193 Hong, Xinying Guo, Lei Yang, and Ziwei Liu. MotionDif-
1194 fuse: Text-driven human motion generation with diffusion
1195 model. *IEEE Transactions on Pattern Analysis and Ma-*
1196 *chine Intelligence*, 2024. 3
- [110] Ziqiang Zhang, Long Zhou, Chengyi Wang, Sanyuan Chen,
1200 Yu Wu, Shujie Liu, Zhuo Chen, Yanqing Liu, Huaming
1201 Wang, Jinyu Li, et al. Speak foreign languages with your
1202 own voice: Cross-lingual neural codec language modeling.
1203 *arXiv e-prints*, pages arXiv–2303, 2023. 2
- [111] Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun,
1204 Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu,
1205 et al. LIMA: Less is more for alignment. *arXiv preprint*
1206 *arXiv:2305.11206*, 2023. 2
- [112] Adrian Łańcucki. Fastpitch: Parallel text-to-speech with
1207 pitch prediction. In *ICASSP 2021 - 2021 IEEE Interna-*
1208 *tional Conference on Acoustics, Speech and Signal Pro-*
1209 *cessing (ICASSP)*, pages 6588–6592, 2021. 2, 5, 6
- 1210
1211