

# Empathic Machines: Using Intermediate Features as Levers to Emulate Emotions in Text-To-Speech Systems

Anonymous ACL submission

## Abstract

We present a method to control the emotional prosody of Text to Speech (TTS) systems by using phoneme-level intermediate features (pitch, energy, and duration) as levers. As a key idea, we propose Differential Scaling (DS) to disentangle features relating to affective prosody from those arising due to acoustics conditions and speaker identity. With thorough experimental studies, we show that the proposed method improves over the prior art in accurately emulating the desired emotions while retaining the naturalness of speech. We extend the traditional evaluation of using individual sentences for a more complete evaluation of HCI systems. We present a novel experimental setup by replacing an actor with a TTS system in offline and live conversations. The emotion to be rendered is either predicted or manually assigned. The results show that the proposed method is strongly preferred over the state-of-the-art TTS system and adds the much-coveted “human touch” in machine dialogue. Audio samples from our experiments are available at: <https://emttts.github.io/tts-demo/>

## 1 Introduction

“The text is like a canoe, and the river on which it sits is the emotion. It all depends on the flow of the river, which is your emotion. The text takes on the character of your emotion.”

— Sanford Meisner

In natural language processing, vocabulary and grammar tend to take center stage, but those elements of speech only tell half the story. Affective prosody provides context and gives meaning to words, and keeps listeners engaged. Understanding emotional prosody is central to language and social development. Studies suggest that we show remarkable sensitivity to prosody “even as infants” (Nazzi et al., 1998; Massicotte-Laforge and Shi, 2015). Recently Kraus (2017) shows that voice-only communication likely elicits higher empathic

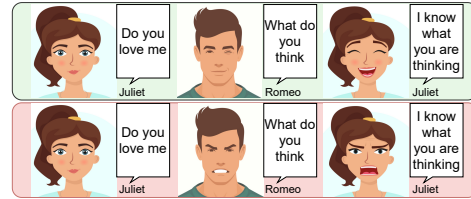


Figure 1: Dialogues can have different meaning despite having same text. Also, starting at same emotion Juliet has different emotion post Romeo’s response.

accuracy than even multi-sense modes including facial expressions.

Buchholz (2016) shows that any meaningful spoken dialogue cannot happen without some amount of prosodic matching. As humans, we naturally anticipate and adapt with emotional cues in conversing with others, see Figure 1 for an example. Celebrated trainer Sanford Meisner employed this to develop *Meisner technique* for theatre actors to react naturally to others in the environment as opposed to *method acting*. The importance of emotional prosody in conversations cannot be overstated and TTS models need to fill this gap to make human-like conversations possible in HCI systems.

Mitchell and Xu (2015) study the value of emotional prosody in HCI and emphasize its role in healthcare dialogue systems, improving social interaction skills in people with autism, augmentative and alternative communication devices and gaming narratives. They explain that successfully incorporating expressive speech into HCI, involves two aspects: (a) prosodic emotion recognition and (b) expression of emotional prosody. Considerable effort has been made towards recognizing and predicting the emotional nuances in human dialogues (Kim and Vossen, 2021; Poria et al., 2019b; Zhu et al., 2021; Li et al., 2017; Poria et al., 2021; Vinyals and Le, 2015). However, current TTS systems are yet to improve on rendering emotive or expressive speech for real-world HCI systems.

State-of-the-art TTS systems (Ren et al., 2020;

Wang et al., 2017) tend to exhibit average emotions for a given phoneme sequence by taking the mean of utterances from training data. Some efforts towards improving expressiveness (like Battenberg et al., 2019; Karlapati et al., 2020) provide prosody control using a reference clip. Others like Sivaprasad et al. (2021) and Habib et al. (2019) further focused on controllability exposing levers that can be manipulated at inference-time to derive the intended expression. However, the quality and stability of synthesized speech heavily depends on various modeling choices. Emotion or prosody modeling, for example, could pick from numerous available discrete or continuous space representations. The encoder network module chosen might vary in its ability to disentangle prosody from other acoustic features like speaker identity and adaptability to content. For example, those relying on reference clip to replicate prosody might perform poorly when input text is unsuitable for rendering with prosody of reference. Some models feed prosody features with phoneme embeddings directly into the decoder while others use them to predict intermediate features that are used in conditioning the decoder. It is empirically verified (like in Sivaprasad et al., 2021) that intermediate features could be suitably manipulated to bring about the desired change in expression.

We take this direction forward to endow the intermediate feature prediction module with affective state control over the final rendering. We propose *Differential Scaling (DS)* of the predicted intermediates to bring about the required change in emotion. The *DS* module is aimed to effect only emotion as intended while remaining agnostic to all other features like speakers identities or acoustic conditions as seen in train data. We show that this significantly improves the naturalness of the generated speech, while allowing finer control over prosody.

In addition to comparing our model’s renderings against various others’ from literature for naturalness and emotion control on conventional single utterances drawn from disconnected contexts, we also evaluate them in conversations. We curate data with conversational theatre dialogues and replace an actor with a TTS system. We use its response as a proxy to evaluate the empathic accuracy. In another experiment, we had a theatre director control the emotion levers of our TTS model in a live conversation with the actor to evaluate controllability. As demonstrated in the results, our proposed

method significantly improves over existing methods in producing suitable prosodic variation lending closer to human-like conversations. The rest of this paper will elaborate on the following contributions of this work.

- We propose a simple technique of using a *DS* module to better emulate emotions in TTS rendered speech. This works as plug-and-play with both autoregressive and non-autoregressive TTS models that predict prosodic features as an intermediate step.
- Our work extends the literature of training controllable and expressive TTS models with improved empathic accuracy and without specific studio recorded data.
- Finally, we present novel methods and data for evaluating TTS models in real conversations with human subjects. The method of evaluation is a useful step towards filling the gap of emulating emotional speech that needs more work.

## 2 Related Work

**Prosody and conversational speech.** Unlike in written text, spoken words contain additional non-verbal information. These cues are collectively termed prosody (Leentjens et al., 1998) that include variations in tone, pitch, energy, duration, accents, intonation, stress, etc. Buchholz (2016) showed that prosodic exchange is unavoidable in human dialogue. Various machine learning methods have been proposed to predict emotion in speech from its prosody variations (Asgari et al., 2014; Kamarudin and Abdul Rahman, 2013). Variations in pitch accents (Nielsen et al., 2020), for example, lead to a significant difference in how the receiver perceives the content. A sentence (like I said **un**lock the door, not lock it from (Rosenberg and Hirschberg, 2009)) could be delivered both as a statement and a command by merely changing prosody.

Emotion recognition in conversations has gained increasing attention for developing empathetic machines with emotion-tagged multi-modal data publicly available for modeling like (Li et al., 2017; Poria et al., 2019a; Busso et al., 2008). While most methods like (Majumder et al., 2019; Jiao et al., 2019) use a combination of text and speech information, some leverage additional side-information

169 from broader context (Ghosal et al., 2020) and the  
170 topic of conversation (Zhu et al., 2021).

171 In such labeled data, emotion is often rep-  
172 resented as a categorical variable over a dis-  
173 crete space following models like Ekman’s ba-  
174 sic emotions (Ekman, 1992) or the wheel of  
175 Plutchik (Plutchik, 1980). This choice is largely  
176 owing to the ease of annotating data. Russell (1980)  
177 proposed a continuous two-dimensional space as  
178 an alternative called valence-arousal model for hu-  
179 man emotions. Arousal signifies the intensity of  
180 the emotion while valence captures its polarity. It  
181 has been extended to add a third dimension of dom-  
182 inance, making it the valence-arousal-dominance  
183 (VAD) model. VAD has since been widely used in  
184 modelling emotion in music (Grekow, 2016; Rach-  
185 man et al., 2019), speech (Asgari et al., 2014; Ka-  
186 maruddin and Abdul Rahman, 2013) and other con-  
187 tent (Joshi et al., 2019; Buechel and Hahn, 2017).  
188 We use the continuous space representation as it is  
189 richer and more convenient to handle in our model.

190 **Expressive and controllable TTS.** Neural  
191 TTS systems are now increasingly popular, im-  
192 proving upon older concatenative statistical sys-  
193 tems (Michelle and Georgia, 2020) in synthesized  
194 speech naturalness. These are broadly sequence-  
195 to-sequence networks with an encoder processing  
196 the input text or phoneme sequence followed by a  
197 decoder that generates the sequence of Mel frames  
198 for output speech. Mel frames are then projected  
199 into the time domain by a vocoder (van den Oord  
200 et al., 2016; Griffin and Lim, 1984) to generate  
201 the speech. Decoding could be autoregressive with  
202 Tacotron-like models (Wang et al., 2017) or non-  
203 autoregressive with FastSpeech-like models (Ren  
204 et al., 2019).

205 Non-autoregressive models are faster at infer-  
206 ence than autoregressive models with about com-  
207 parable naturalness of speech quality (Ren et al.,  
208 2020). The trick non-autoregressive models use  
209 to generate Mel frames in parallel is to predict the  
210 relevant features as an intermediate step and con-  
211 dition the independent decoding of Mels on them.  
212 This technique is now increasingly adopted for au-  
213 toregressive models as well (Wang et al., 2021)  
214 to predict features like phoneme duration that im-  
215 prove decoding stability avoiding alignment issues.  
216 Our method is compatible with any architecture  
217 that predicts prosodic features of pitch, energy, and  
218 duration as an intermediate step before decoding.

219 Going beyond the naturalness of speech, there

220 has been considerable effort to improve the expres-  
221 siveness of the renderings. Some focused on learn-  
222 ing a linear space of variations in speech expres-  
223 sions for selecting a suitable variation at inference  
224 time. Wang et al. (2018) learn this space unsuper-  
225 vised by encouraging it to explain all variations in  
226 training data not captured in content embedding.  
227 A reference encoder maps an input utterance to a  
228 style embedding as a linear combination of basis  
229 style vectors. Manual analysis is required to un-  
230 derstand the prosody feature learned into a basis  
231 vector that could include variations like vocal depth  
232 or pitch, speaking rate, or even background noise  
233 as available in training data. While this offers style  
234 control, it does not explicitly learn the prosody vari-  
235 ations of interest into the style space. Our work  
236 focuses on the same level of control but specifically  
237 over the affective state as labeled in some data for  
238 supervision.

239 Sivaprasad et al. (2021) propose a model sim-  
240 ilar to (Wang et al., 2018) with style tokens re-  
241 stricted to valence and arousal. However, the ab-  
242 solute (pitch, energy, duration) feature predictions  
243 restrict prosody control, leading to unnatural dis-  
244 tortions. Specifically, it skews more towards retain-  
245 ing the speaker’s voice identity than the emotion  
246 and entangles emotion with other acoustic features.  
247 Karlapati et al. (2020) replace the linear style space  
248 with a variational reference encoder to generate  
249 prosody embedding to condition the decoder. Bat-  
250 tenberg et al. (2019) use a similar variational model  
251 but instead force its posterior to match that of the  
252 reference utterance to copy prosody with a con-  
253 trollable parameter determining the closeness of  
254 the match. This trick alleviates certain issues like  
255 in pitch-range (Younggun and Taesu, 2019) and  
256 transfer to unrelated sentences but exposes a lower  
257 degree of control with no explicit levers to operate,  
258 as possible in our work.

259 Habib et al. (2019) propose to learn explicit la-  
260 tent representation for various prosodic variables,  
261 segregating them into explicitly controllable (like  
262 affect, speaking rate, etc) and implicit (like intona-  
263 tion, rhythm, stress, etc). While the model offers a  
264 higher degree of explicit control, it requires using  
265 a proprietary studio recorded data with utterances  
266 reflecting prompted emotions at specified arousal.  
267 Dependence on explicit supervision from studio  
268 recorded data makes it harder scale this model  
269 across languages and other prosodic variations. In  
270 contrast, we use publicly available data with emo-

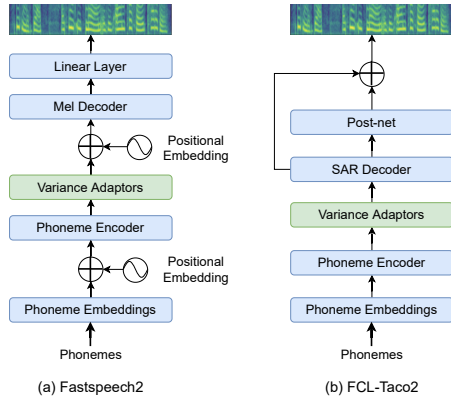


Figure 2: Backbone TTS architectures.

tion labels to train our models.

There are other methods that try to predict suitable prosody features from text content. Raitio et al. (2020) add a prosody encoder module to standard TTS network that predicts certain hand-crafted prosody features from text embedding of input. This prosody encoder is used with a small optional bias for affect variations at inference. Hodari et al. (2021) extend this to replace hand-crafting prosody features with explicit training followed by their prediction from text. Karlapati et al. (2021) further enrich the textual context using BERT embeddings and parse-trees. These methods are limited in expressiveness offering no control over rendering emotion that our work focuses on.

### 3 Model

Our network uses a backbone TTS that can be borrowed from any model which predicts pitch, energy and duration as intermediates features from input phoneme sequence. This network learns to predict the average features for given phonemes. Following the convention in earlier works, we refer to the intermediate features as variances and the module that predicts them as variance adaptor. Prior work improves standard variance adaptors in, say FastSpeech2, by conditioning on emotion variables of valence-arousal in addition to the phoneme sequence to generate expressive speech. We refer to it as Emotional Variance Adaptor (EVA) for which we propose an alternative. Our proposed Differential Scaler (DS) module determines how best to vary the output of the EVA to bring the desired change in emotion. We describe the details of these network choices in this section; specifically, the broader backbone network architecture and the different variance adaptor modules from non-emotive

baseline, emotive baseline and our proposal.

#### 3.1 Backbone

We present experiments with two suitable choices for our backbone systems, FastSpeech2 and FCL-Taco2. The backbone has three modules; an encoder, variance adaptor and decoder. The encoder maps an input phoneme sequence to its embedding. Given this representation, the variance adaptor predicts the pitch, energy and duration for each of the phonemes. These intermediate features are processed by the decoder module downstream to return Mel-spectrogram frames. We reuse the encoder and decoder modules as designed in their original architectures without any changes. We refer readers to the respective papers for details of these networks. Wavenet (van den Oord et al., 2016) vocoder is used to map Mel-spectrogram outputs of the decoder to time-domain raw audio.

#### 3.2 Variance adaptor module

**Non-emotive baselines.** Our baseline models of FastSpeech2 and FCL-Taco2 are trained with the variance adaptors as described by their authors. We also train a derivative of the FastSpeech2 with the variance adaptor modified to make predictions at the phoneme-level and not at frame-level. A duration  $d_\pi$  is predicted for each phoneme  $\pi$ , following which the length regulator repeats the hidden state of that phoneme  $\pi$  times. Also unlike FastSpeech2, we use this length regulator after the predicted pitch and energy are added to the encoder output. We refer to this derivative as FastSpeech2 $\pi$ .

**Emotive baseline.** Sivaprasad et al. (2021) conditioned the variance adaptor of FastSpeech2 on additional emotion embedding that gives the model control over prosody of the rendered speech. It generates the emotion embedding as a linear weighted combination of the the valence and arousal vectors that are learned from data during training. The weights are valence and arousal values as annotated for training and can be used as control levers to modify emotion during inference. This emotion variance adaptor (EVA) module generates suitable intermediate features of pitch and energy at frame-level and duration at phoneme-level. These features are consumed by the decoder along with the encoder output in generating Mel frames. While this helps control emotional prosody rendered speech, it leads to a significant drop in perceptual quality and naturalness relative to the baselines. Our contribution is an alternative design of the variance

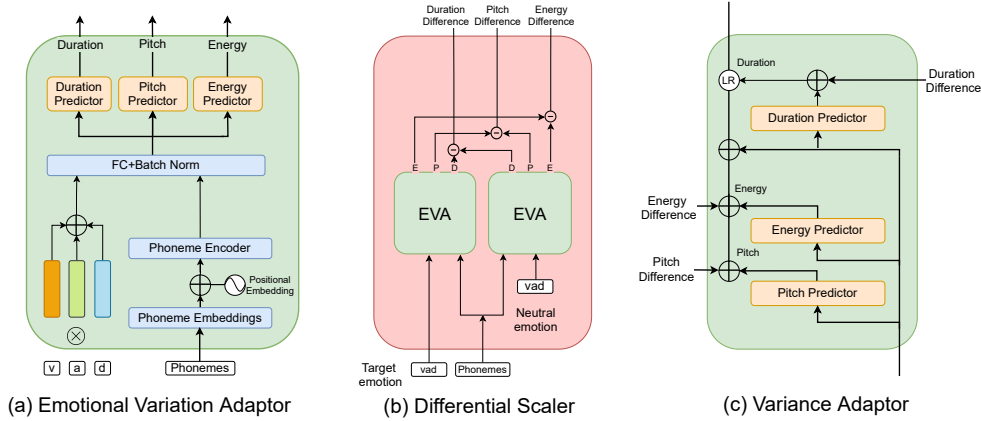


Figure 3: Schematic diagram of the proposed model.

357 adaptor module that improves upon [Sivaprasad et al.](#)  
 358 (2021)’s FastSpeech2 + EVA model in emotion control and expressiveness and upon the baselines in  
 359 terms of naturalness.  
 360

361 **Differential Scaler.** We extend the emotion representation from EVA to include dominance in addition to valence and arousal values. Dominance is the degree of control exerted by an emotion. Including dominance dimension to the emotion space expands the range of emotions the TTS model can express. For example, by introducing this dimension, we can better distinguish outputs for emotions like ‘anger and fear’ or ‘sad and contempt’.

370 The *Differential Scaler* module further extends EVA to estimate the change in variances necessary for a pronounced effect of the target emotion relative to its neutral counterpart. As shown in Figure 3(b), the variances are estimated using the EVA module for a given phoneme sequence at two different triplets of VAD values. One prediction corresponds to the neutral emotion with VAD values all set to zeros. The other prediction corresponds to the chosen VAD values of target emotion. We take the difference of these two estimates as the direction along which the variances can be varied for the desired change in emotion without effecting other acoustic features. We are implicitly making two assumptions here. Emotion variations are captured as linear transformations in this space and that there is a strong disentangling of emotional prosody with other acoustic features in this space. Results from our empirical evaluation favorably support the above assumptions.

## 390 4 Training

391 Modelling with intermediate features facilitates  
 392 training the backbone and the variance adaptors

393 independently on different data. We exploit this  
 394 to train our variance adaptor on scarcely available  
 395 VAD annotated data while reusing backbone models  
 396 trained on abundant transcribed speech data.

397 **Backbone.** We train two backbone networks  
 398 FastSpeech2 $\pi$  (non-autoregressive) and FCL-Taco2  
 399 (autoregressive) on Blizzard 2013 dataset ([King and Karaiskos, 2014](#)). It contains 147 hours of  
 400 Catherine Bayers’s speech, reading books in American English. Due to the style of reading, the dataset  
 401 is rich in expressiveness and spans different combinations of pitch, energy and duration. Both models  
 402 are trained with Mel loss (mean absolute error between predicted and ground truth Mels), pitch loss,  
 403 energy loss and duration loss (mean square error between predicted and ground truth features). Both  
 404 models are trained for 200K iterations using Adam optimizer with warm-up learning rate scheduler  
 405 and batch size of 16.  
 406

407 **EVA.** We train EVA on MSP-Podcast corpus  
 408 ([Lotfian and Busso, 2019](#)) annotated with arousal, valence and dominance values. The corpus consists  
 409 of around 100 hours speech data but their transcriptions are not available. We generate transcripts  
 410 using a speech-to-text model. We use Montreal-Forced-Aligner (MFA) ([McAuliffe and Sonderegger, 2017](#))  
 411 for phoneme alignments. Those transcripts that MFA fails to find a good alignment for are filtered out.  
 412 The remaining utterances add up to about 71 hours of emotive speech data which we use to train our  
 413 EVA. We train pitch, energy and duration predictors conditioned on VAD values minimizing only the sum  
 414 of variance losses. For all the experiments, text transcripts are converted to phonemes using  
 415 ([Sun et al., 2019](#)). We generate Mel spectrogram from the audio files similar to  
 416 ([Wang et al., 2017](#)). Pitch and energy are computed  
 417  
 418  
 419  
 420  
 421  
 422  
 423  
 424  
 425  
 426  
 427  
 428  
 429

430 from the Mel spectrogram and we use MFA for  
431 aligning phonemes to train the duration predictor.

## 432 5 Experiments and user study

433 We present three experiments; comparison with  
434 prior-art using conventional evaluation metrics,  
435 those for emotional consistency with pre-recorded  
436 audio, and finally, live conversations with humans.

### 437 5.1 Comparisons with prior-art

438 We compare the proposed approach against four  
439 state of the art TTS models. The list includes two  
440 non-emotive TTS models (FastSpeech2 and FCL-  
441 Taco2), one reference-based method (Cai et al.,  
442 2021) and one AV conditioned model (FastSpeech2  
443 + EVA). We also compare our method with the  
444 modified backbone, Fastspeech2 $\pi$ .

445 To evaluate the perceptual-quality/naturalness  
446 we compare Mean Opinion Score (MOS) (Chu and  
447 Peng, 2006) averaged across forty subjects profi-  
448 cient in English. We synthesize twenty different  
449 sentences from the test set using each of the seven  
450 models. We prepare user study by picking five  
451 samples rendered by each model to make a survey.  
452 Annotator rates each sample on a Likert scale of  
453 one for ‘completely unnatural’ to five for ‘com-  
454 pletely natural’.

455 To evaluate the emotional expressiveness of the  
456 proposed model, we perform two surveys. In the  
457 first survey, given a sample, we ask the user to  
458 choose the best perceived emotion from a set of  
459 four, namely, ‘Happy’, ‘Sad’, ‘Angry’ and ‘Fear’.  
460 We ask the raters to not judge the textual content  
461 and annotate the emotion for each sample based on  
462 the rendering alone. In the second survey we evalu-  
463 ate the efficacy of the models to bring about finer  
464 control over emotion. We generate two samples  
465 with same broader emotion category but with two  
466 levels of intensity. The subject now has to identify  
467 the sample with higher intensity. For both surveys  
468 we generate five samples per emotion and twenty  
469 samples for each model. We aggregate the rating  
470 across forty proficient English language speakers.

### 471 5.2 Emotional consistency in dialogues

472 Previous efforts in prosody controlled TTS have  
473 been evaluated on individual sentences without con-  
474 text. We propose a novel evaluation strategy us-  
475 ing excerpts from theater recordings. We replace  
476 the audio of one of the actors in the conversation  
477 with renderings from a TTS model and have a hu-

478 man subject evaluate it for emotional consistency.  
479 The emotion for TTS renderings are chosen manu-  
480 ally by a theater director. We compare this with  
481 TTS rendered with emotion predicted using Tod-  
482 Kat (Zhu et al., 2021) from the dialogues spoken  
483 so far. This study consolidates the two aspects of  
484 HCI we mentioned in the introduction; prosodic  
485 emotion recognition and its expression in TTS ut-  
486 terances.

487 The dataset is curated using segments from four  
488 popular plays, namely, ‘Speed-the-Plow’, ‘Night,  
489 Mother’, ‘Bobby Gould in Hell’ and ‘Death of a  
490 Salesman’. We select 30 dialogue segments collec-  
491 tively from the four plays with an average dialogue  
492 length of 90 seconds per segment. Timestamps of  
493 segments selected from each play is given in sup-  
494plementary material. We replace the female voice  
495 in the segment with (a) non-emotive TTS model  
496 (Fastspeech2 $\pi$ ) (b) our model with emotion pre-  
497 dicted for each utterance using TodKat and (c) our  
498 model with a senior theatre director picking the  
499 emotion for each utterance. We randomly pick five  
500 dialogues from the 30 samples in all three settings  
501 for each of our surveys. We ask forty raters to rank  
502 the three setting in terms of the emotional consis-  
503 tency of the dialogue *i.e.*, to judge the naturalness  
504 and aptness of the emotional prosody in the given  
505 context.

### 506 5.3 Conversation with Meisner trained actor

507 A Meisner trained actor responds to another actor  
508 taking into account his/her behavior. In this experi-  
509 ment, we observe how a Meisner trained actor (Ac-  
510 tor M) reacts in a live dialogue initiated by (a) an-  
511 other trained human actor, (b) a non-emotive TTS  
512 (Fastspeech2 $\pi$ ) and (c) our model (Fastspeech2 $\pi$   
513 with DS). We use the same neutral script with 18  
514 lines in all three cases. We use the behavior of  
515 Actor M during interaction with the human as refer-  
516 ence. The closeness of Actor M’s behavior to this  
517 reference while interacting with the two TTS mod-  
518 els is used as a measure of the latter’s effectiveness  
519 in rendering speech expressive enough to evoke an  
520 emotive response.

521 For each of the three scenarios, the conversation  
522 is initiated with two different emotional states, viz.  
523 (a) highly positive and (b) highly negative. The  
524 emotion for our TTS model is chosen live on-the-  
525 fly by a theatre director from fourteen bins in the  
526 discretized arousal-valence space. The bins are  
527 chosen to span the V-shape around high-arousal-

Model	MOS	Finer Control	Coarse Control				Average
			Happy	Sad	Angry	Fear	
Fastspeech2	3.80	-	-	-	-	-	-
FCL-Taco2	3.39	-	-	-	-	-	-
Fastspeech2 $\pi$	3.84	-	-	-	-	-	-
Cai et al., 2021	3.08	80.0	22.7	40.9	52.3	-	38.7
FastSpeech2 + EVA	3.01	81.2	20.0	68.7	52.9	-	47.2
Fcl-Taco2 + DS (our model)	3.30	83.5	90.1	53.3	56.5	46.8	61.8
Fastspeech2 $\pi$ + DS (our model)	<b>3.91</b>	<b>85.0</b>	68.4	50.0	59.5	79.1	<b>64.2</b>

Table 1: Results for qualitative analysis comparing our model with prior art.

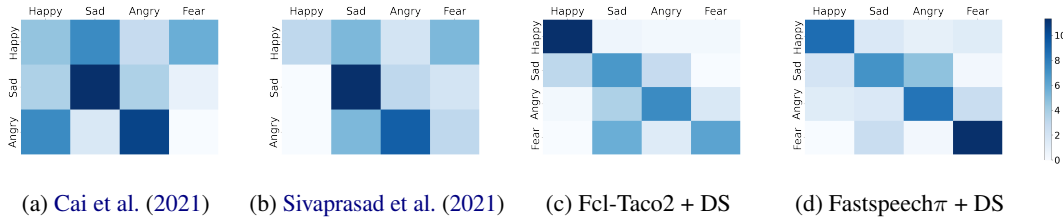


Figure 4: Confusion matrices of models performance in the survey to pick the correct emotion. Rows are true emotions and columns are picked emotions. Figure to be viewed in color.

high-valence and low-arousal-neutral-valence (Dietz and Lang, 1999). We take majority vote of three listener ratings for each utterance of Actor M on the same discretized arousal-valence space to allow quantitative comparisons.

## 6 Results

### 6.1 Comparing with prior art

**Naturalness.** Table 1 compares the audio quality of the TTS models listed in Section 5.1. It can be seen that the proposed model achieves affective control, without drop in perceived audio quality. In contrast, previous SOTA emotive models (Cai et al. (2021) and Fastspeech2 + EVA) achieve control over emotion at the cost of naturalness (MOS of 3.08 and 3.01 respectively). This result demonstrates the efficacy of using DS module over EVA and validates its ability to disentangle affective features from the acoustic ones. The MOS score of Fastspeech2 $\pi$  improves with addition of DS, as some samples appear more natural when rendered in intended emotions.

**Coarse affective control.** Results corresponding to emotion detection are presented in Table 1. For each sample, the raters were asked to choose one among the four discrete emotions. On an average, the Fastspeech2 $\pi$  + DS gives best results, outperforming the other models by a significant margin. We observe 17%, and 25.5% improvement over (Cai et al. (2021), Fastspeech2 + EVA) respectively. Figure 4 shows the confusion matrix for this

survey. Our models are better at differentiating positive valence emotions from the negative ones. There is still a scope of improvement in distinctly expressing low valence emotions.

**Finer affective control.** When asked raters to pick the sample from a pair that expresses a particular emotion better, 85% of the times they were able to pick the sample that was actually rendered with a higher arousal value (Table 1). Our best performing model shows 3.8% improvement over Fastspeech2 + EVA and 4.0% improvement over Cai et al. (2021).

### 6.2 Emotional consistency in dialogues

As described in Section 5.2, we evaluate the emotional consistency of a dialogue when a TTS model replaces an actor in excerpts from a play. Figure 5 shows that emotive models bring significant improvement in emphatic quality of conversations and are picked 80% of the times as the first preference. This result reiterates the hypothesis (Wang et al., 2018) that prosody averaging as in non-emotive TTS is insufficient for emulating emotionally consistent conversations.

Another important observation is how emphatic quality measured as user’s first preference falls from 52% to 27% in moving away from hand-picked to model-predicted emotions. This suggests a scope for improvement for emotion prediction models. Nonetheless the results present clear evidence that tying together emotion prediction models to expressive TTS is significantly more prefer-

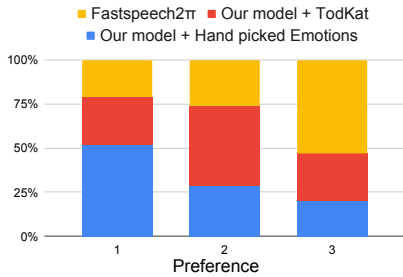


Figure 5: Comparison of emotional consistency in conversations across the three settings described in section 5.2. Figure best viewed in color.

able to a non-emotive TTS.

This proposed evaluation methodology is more comprehensive and enables assessment of a consolidated conversational system as required in expressive HCI that includes various moving parts like causal emotion recognition in conversation and expressive TTS. This is not feasible with the traditional approach of evaluating on individual sentences drawn from distinct contexts. We argue that this evaluation with contextual dialogues from a conversation is more coherent to humans as reflected in inter-annotator agreement measured by Fleiss’s Kappa Score (FKS). FKS goes up by 34% from 0.43 in traditional coarse affective control (Table 1) to 0.58 for our evaluation strategy (Figure 5). We hope this will be useful in a more thorough evaluation of expressive HCI systems.

### 6.3 Conversation with Meisner trained actor

As mentioned in Section 5.3, we gather the behavioural response of a Meisner trained human actor to TTS systems (emotive and non-emotive) and compare it against his/her reference response to another human actor. We use Pearson’s correlation  $\rho$  with reference for valence and compare mean-std ( $\mu, \sigma$ ) for arousal values.

When the conversation was triggered with a positive initial emotion, we had a high  $\rho(\text{FastSpeech2}\pi+\text{DS}, \text{human})$  of 0.702 for our model compared to negative correlation for non-emotive TTS at  $\rho(\text{FastSpeech2}\pi, \text{human})$  of  $-0.282$ . Similarly for a negative initial emotion  $\rho(\text{FastSpeech2}\pi+\text{DS}, \text{human})$  was high 0.838 relative to low  $\rho(\text{FastSpeech2}\pi, \text{human})$  of 0.158.

We find that the average arousal for the human response to our TTS ( $\mu=3.5, \sigma=1.06$ ) is comparable to a human-human conversation ( $\mu=3.94, \sigma=0.97$ ), as opposed to the response to a non-emotive TTS ( $\mu=2.55, \sigma=0.49$ ). This indicates that the range of

arousal response elicited from a human actor by our TTS is comparable to a human-human conversation as opposed to that of a prosody unaware TTS.

We also interviewed the human actor about the experience of conversing with the TTS systems. He reported that our TTS gave him "an emotional structure". He felt that the TTS could "dictate the neutral part of the script to change it". He could "remember specific utterances" by our TTS and their emotional content which "drove him" to respond in an emotional manner. In contrast, he reported that the prosody unaware TTS gave "dry answers", made him feel that it was "disinterested", "auto generated" and "did not evoke excitement". He expressed that he "could not have a longer conversation with it".

## 7 Conclusion

This work presents a novel way to lever the prosodic features (pitch, energy and duration) to modify emotions in the output of a TTS system. Our method is model agnostic and can be used with any TTS backbone that predicts prosodic features in an intermediate step. This method outperforms existing approaches by a significant margin in its ability to accurately render desired emotions, while preserving the naturalness of speech. We curated theatre conversation data to evaluate and show that our prosody-aware TTS better maintains the natural flow of emotions in conversations. Our work shows promise in consolidation of prosodic emotional recognition and expression, a coveted pursuit in the field of HCI. We present further qualitative experiments involving professional theatre artists and demonstrate that the proposed TTS method leads to more human-like conversations. While exposing valence, arousal and dominance values as model levers improves control over the final rendering, in reality it is overwhelming for the user to choose them correctly for a desired output. This is further aggravated by the fact that some sentences cannot be suitably spoken with a chosen set of values, degrading output quality. These are limitations that need to be addressed and appropriately deriving these values from semantics of text input or reference clips could be relevant future directions. Affective control is incomplete without explicit levers on the intonations, which is another limitation to be looked upon in the future work.



## 8 Ethical concerns

This work shares the same concerns as with others in the domain of TTS systems as discussed by Habib et al. (2019). With TTS outputs getting closer to actual human speech, there could be a potential misuse. The threat of abuse of fake voices is particularly high with similar developments in conjugate areas like computer vision. However, the benefits of improvements to emotive TTS technology could significantly benefit HCI and the corresponding applications to problems in healthcare and other domains. Example applications include healthcare dialogue systems, improving social interaction skills in people with autism and augmentative communication devices. TTS systems synthesizing speech with empathy can ease machine interaction in many touchpoint applications. While the benefits seem to outweigh the concerns at this point, we believe the research community should proactively continue to identify methods for detection and prevention of misuse.

## References

Meysam Asgari, Géza Kiss, Jan Van Santen, Izhak Shafran, and Xubo Song. 2014. Automatic measurement of affective valence and arousal in speech. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 965–969. IEEE.

Eric Battenberg, Soroosh Mariooryad, Daisy Stanton, RJ Skerry-Ryan, Matt Shannon, David Kao, and Tom Bagby. 2019. Effective use of variational embedding capacity in expressive end-to-end speech synthesis. *arXiv preprint arXiv:1906.03402*.

Michael B Buchholz. 2016. Conversational errors and common ground activities in psychotherapy—insights from conversation analysis. *International Journal of Psychological Studies*, 8(3):134–153.

Sven Buechel and Udo Hahn. 2017. Emobank: Studying the impact of annotation perspective and representation format on dimensional emotion analysis. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 578–585.

Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeanette N Chang, Sungbok Lee, and Shrikanth S Narayanan. 2008. Iemocap: Interactive emotional dyadic motion capture database. *Language resources and evaluation*, 42(4):335–359.

Xiong Cai, Dongyang Dai, Zhiyong Wu, Xiang Li, Jingbei Li, and Helen M. Meng. 2021. Emotion controllable speech synthesis using emotion-unlabeled

dataset with the assistance of cross-domain speech emotion recognition. *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5734–5738.

Min Chu and Hu Peng. 2006. Objective measure for estimating mean opinion score of synthesized speech. US Patent 7,024,362.

Richard Dietz and Annie Lang. 1999. Affective agents: Effects of agent affect on arousal, attention, liking and learning. In *Proceedings of the Third International Cognitive Technology Conference, San Francisco*.

Paul Ekman. 1992. An argument for basic emotions. *Cognition & emotion*, 6(3-4):169–200.

Deepanway Ghosal, Navonil Majumder, Alexander Gelbukh, Rada Mihalcea, and Soujanya Poria. 2020. Cosmic: Commonsense knowledge for emotion identification in conversations. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 2470–2481.

Jacek Grekow. 2016. Music emotion maps in arousal-valence space. In *IFIP International Conference on Computer Information Systems and Industrial Management*, pages 697–706. Springer.

Daniel Griffin and Jae Lim. 1984. Signal estimation from modified short-time fourier transform. *IEEE Transactions on acoustics, speech, and signal processing*, 32(2):236–243.

Raza Habib, Soroosh Mariooryad, Matt Shannon, Eric Battenberg, RJ Skerry-Ryan, Daisy Stanton, David Kao, and Tom Bagby. 2019. Semi-supervised generative modeling for controllable speech synthesis. In *International Conference on Learning Representations*.

Zack Hodari, Alexis Moinet, Sri Karlapati, Jaime Lorenzo-Trueba, Thomas Merritt, Arnaud Joly, Ammar Abbas, Penny Karanasou, and Thomas Drugman. 2021. Camp: a two-stage approach to modelling prosody in context. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6578–6582. IEEE.

Wenxiang Jiao, Haiqin Yang, Irwin King, and Michael R Lyu. 2019. Higr: Hierarchical gated recurrent units for utterance-level emotion recognition. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 397–406.

Tanmayee Joshi, Sarath Sivaprasad, and Niranjana Pedanekar. 2019. Partners in crime: Utilizing arousal-valence relationship for continuous prediction of valence in movies. In *AffCon@ AAAI*.

780	Norhaslinda Kamaruddin and Abdul Wahab Abdul Rahman. 2013. <a href="#">Valence-arousal approach for speech emotion recognition system</a> . In <i>2013 International Conference on Electronics, Computer and Computation (ICECCO)</i> , pages 184–187.	833
781		834
782		835
783		836
784		837
		838
785	Sri Karlapati, Ammar Abbas, Zack Hodari, Alexis Moinet, Arnaud Joly, Penny Karanasou, and Thomas Drugman. 2021. Prosodic representation learning and contextual sampling for neural text-to-speech. In <i>ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 6573–6577. IEEE.	839
786		840
787		841
788		842
789		
790		
791		
792	Sri Karlapati, Alexis Moinet, Arnaud Joly, Viacheslav Klimkov, Daniel Sáez-Trigueros, and Thomas Drugman. 2020. <a href="#">CopyCat: Many-to-Many Fine-Grained Prosody Transfer for Neural Text-to-Speech</a> . In <i>Proc. Interspeech 2020</i> , pages 4387–4391.	843
793		844
794		845
795		846
796		847
797	Taewoon Kim and Piek Vossen. 2021. Emoberta: Speaker-aware emotion recognition in conversation with roberta. <i>arXiv preprint arXiv:2108.12009</i> .	848
798		849
799		850
800	Simon King and Vasilis Karaiskos. 2014. The blizzard challenge 2013.	851
801		852
802	Michael W Kraus. 2017. Voice-only communication enhances empathic accuracy. <i>American Psychologist</i> , 72(7):644.	853
803		854
804		855
805	Albert FG Leentjens, Sandra M Wielaert, Frans van Harskamp, and Frederik W Wilmink. 1998. Disturbances of affective prosody in patients with schizophrenia; a cross sectional study. <i>Journal of Neurology, Neurosurgery &amp; Psychiatry</i> , 64(3):375–378.	856
806		857
807		858
808		859
809		860
810		861
811	Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. Dailydialog: A manually labelled multi-turn dialogue dataset. In <i>Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 986–995.	862
812		863
813		864
814		865
815		866
816		867
817	R. Lotfian and C. Busso. 2019. <a href="#">Building naturalistic emotionally balanced speech corpus by retrieving emotional speech from existing podcast recordings</a> . <i>IEEE Transactions on Affective Computing</i> , 10(4):471–483.	868
818		869
819		870
820		871
821		872
822	Navonil Majumder, Soujanya Poria, Devamanyu Hazarika, Rada Mihalcea, Alexander Gelbukh, and Erik Cambria. 2019. Dialoguernn: An attentive rnn for emotion detection in conversations. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 33, pages 6818–6825.	873
823		874
824		875
825		876
826		877
827		878
828	Sarah Massicotte-Laforge and Rushen Shi. 2015. The role of prosody in infants’ early syntactic analysis and grammatical categorization. <i>The Journal of the Acoustical Society of America</i> , 138(4):EL441–EL446.	879
829		880
830		881
831		882
832		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

887	Yi Ren, Chenxu Hu, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. 2020. Fastspeech 2: Fast and high-quality end-to-end text to speech. In <i>International Conference on Learning Representations</i> .	Lixing Zhu, Gabriele Pergola, Lin Gui, Deyu Zhou, and Yulan He. 2021. <a href="#">Topic-driven and knowledge-aware transformer for dialogue emotion detection</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 1571–1582, Online. Association for Computational Linguistics.	939
888			940
889			941
890			942
891			943
892	Yi Ren, Yangjun Ruan, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. 2019. Fastspeech: fast, robust and controllable text to speech. In <i>Proceedings of the 33rd International Conference on Neural Information Processing Systems</i> , pages 3171–3180.		944
893			945
894			946
895			947
896			
897	Andrew Rosenberg and Julia Bell Hirschberg. 2009. Detecting pitch accents at the word, syllable and vowel level.		
898			
899			
900	James A Russell. 1980. A circumplex model of affect. <i>Journal of personality and social psychology</i> , 39(6):1161.		
901			
902			
903	Sarath Sivaprasad, Saiteja Kosgi, and Vineet Gandhi. 2021. Emotional prosody control for speech generation.		
904			
905			
906	Hao Sun, Xu Tan, Jun-Wei Gan, Hongzhi Liu, Sheng Zhao, Tao Qin, and Tie-Yan Liu. 2019. <a href="#">Token-Level Ensemble Distillation for Grapheme-to-Phoneme Conversion</a> . In <i>Proc. Interspeech 2019</i> , pages 2115–2119.		
907			
908			
909			
910			
911	Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alexander Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. 2016. <a href="#">Wavenet: A generative model for raw audio</a> . In <i>Arxiv</i> .		
912			
913			
914			
915			
916	Oriol Vinyals and Quoc Le. 2015. A neural conversational model. <i>arXiv preprint arXiv:1506.05869</i> .		
917			
918	Disong Wang, Liqun Deng, Yang Zhang, Nianzu Zheng, Yu Ting Yeung, Xiao Chen, Xunying Liu, and Helen Meng. 2021. <a href="#">Fcl-taco2: Towards fast, controllable and lightweight text-to-speech synthesis</a> . In <i>ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 5714–5718.		
919			
920			
921			
922			
923			
924			
925	Yuxuan Wang, RJ Skerry-Ryan, Daisy Stanton, Yonghui Wu, Ron J Weiss, Navdeep Jaitly, Zongheng Yang, Ying Xiao, Zhifeng Chen, Samy Bengio, et al. 2017. Tacotron: Towards end-to-end speech synthesis. <i>arXiv preprint arXiv:1703.10135</i> .		
926			
927			
928			
929			
930	Yuxuan Wang, Daisy Stanton, Yu Zhang, RJ-Skerry Ryan, Eric Battenberg, Joel Shor, Ying Xiao, Ye Jia, Fei Ren, and Rif A Saurous. 2018. Style tokens: Unsupervised style modeling, control and transfer in end-to-end speech synthesis. In <i>ICML</i> .		
931			
932			
933			
934			
935	Lee Younggun and Kim Taesu. 2019. Robust and fine-grained prosody control of end-to-end speech synthesis. In <i>International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> . IEEE.		
936			
937			
938			