



Prosodic Structure Beyond Lexical Content: A Study of Self-Supervised Learning

Sarenne Wallbridge*, Christoph Minixhofer*, Catherine Lai, Peter Bell

Centre for Speech Technology Research, University of Edinburgh, UK

{s.wallbridge, christoph.minixhofer, peter.bell, c.lai}@ed.ac.uk

Abstract

People exploit the predictability of lexical structures during text comprehension. Though predictable structure is also present in speech, the degree to which prosody—e.g., intonation, tempo, and loudness—contributes to such structure independently of the lexical content is unclear. This study leverages self-supervised learning (SSL) to examine the temporal granularity of structures in the acoustic correlates of prosody. Representations from our proposed Masked Prosody Model can predict perceptual labels dependent on local information, such as word boundaries, but provide the most value for labels involving longer-term structures, like emotion recognition. Probing experiments across various perceptual labels show strong relative gains over untransformed pitch, energy and voice activity features. Our results reveal the importance of SSL training objective timescale and highlight the value of complex SSL-encoded structures compared to more constrained classical structures.

Index Terms: prosody, self-supervised learning, emotion recognition, prominence prediction, phonetic segmentation

1. Introduction

Cognitive science theories often describe human language comprehension as a process of leveraging the predictable structure of linguistic signals, i.e., their *systematicity*, to generate expectations about the upcoming signal [1]. Though empirical support stems primarily from written language comprehension, listeners also exploit systematicity to predict features of upcoming speech, e.g., informational content [2, 3], the position of prosodic focus [4], and the length of upcoming units (sentences [5] and conversational turns [6]). However, the extent to which this predictive capacity is contingent on the structure of lexical information or of its acoustic realisation is unclear. To better understand how listeners generate expectations about upcoming speech, we investigate systematicity in non-lexical features of prosody—relative pitch (F0) and loudness (energy)¹.

We address two challenges in studying the structure of prosody. First, prosody operates at multiple temporal scales, making it difficult to define a task-agnostic unit (cf. [7]). Second, prosody functions *in conjunction* with lexical content—e.g., interpretations of prosodic changes to discourse markers are restricted by lexical semantics [8]. Still, prosody also conveys information *independently* of lexical content, indicating that its acoustic correlates exhibit systematicity. For example, people can be primed to different syntactic disambiguation strategies by prosodic differences in de-lexicalised audio [9].

* for joint authorship.

¹In this paper, we now use ‘prosody’ to refer to non-lexical aspects of speech including relative pitch (F0), energy, and timing features.

Moreover, pitch contours can help classify reported speech [10]. To address both challenges, we apply Self-Supervised Learning (SSL), a mechanism that exploits complex structures across timescales, to study the acoustic correlates of prosody.

We introduce a novel Masked Prosody Model (MPM) which encodes correlates of pitch, loudness and voice activity by learning to reconstruct corrupted feature sequences. By altering the corruption strategy, we ask 1) *Can SSL capture predictable structure in prosody, independent of lexical content?* and 2) *How does corruption timescale affect the utility of representations for predicting perceptual speech labels?* We select labels that rely on structures at different temporal granularities. Using linear probes, we find a task-dependent effect of corruption granularity on the utility of resulting representations and propose a masking strategy for generalisable representations. Comparing MPM to a constrained prosodic encoding reveals that the complex contextualisation of SSL is particularly useful for predicting abstract perceptual labels. Stronger probes illustrate how SSL and task-specific structures interact, underscoring the complexity of prosodic structures. Finally, comparing MPM to SSL representations of the full speech signal highlights how prosodic structure contributes to different tasks.

2. Background

2.1. Self-supervised learning

SSL mechanisms encode the structures of their training data by learning to reverse a corruption function—e.g., removing future, past, or intermediate context. Without a need for external labels, SSL can exploit large amounts of unlabeled data. The resulting representations are useful for downstream tasks with limited labelled data, either as input to task-specific models or through fine-tuning [11]. Transformers trained with SSL capture both global and local structures and have been successfully applied to text [12], speech and a host of other signals [13]; popular models of speech include wav2vec [14] and HuBERT [15]. However, it is unclear which input structures SSL exploits during training. In particular, it is unclear how much prosodic structure is maintained in SSL representations of speech [16].

2.2. Representations of prosody

Neural representations of prosody have recently gained popularity, often with the aim of controlling prosodic realisations of text inputs. Many of these encodings use a subtractive definition of prosody: the variation remaining once phonetic, speaker, and channel information are removed [17]. Disentanglement of these factors has been induced through carefully-tuned bottlenecks or training data design [18, 19]. While such representations are useful for speech generation and tasks like emotion

recognition, the degree of achieved (and achievable) disentanglement is unclear as they show sensitivity to e.g., speaker and lexical perturbation [17, 20]. As such, studying prosody isolated from lexical content with these representations is difficult.

Representations based instead on acoustic correlates of prosody are independent of lexical information but less well-studied. One of the few widely used representations of this type is Continuous Wavelet Transform (CWT) across F0, energy, and duration features for automatic annotation of local prosodic events [21]. CWT captures hierarchical structure by convolving the original signal with sets of wavelets. Using a rule-based labelling method, CWT features achieve comparable performance to supervised methods for word-level prominence and boundary detection and are useful for generating prosodic renditions that are more faithful to references [22]. These results support the importance of hierarchical structure in prosody, which has long been highlighted by linguistic theories [23, 24]. Other work has used the predictability of word-level pitch and intensity features to detect prominent words [25]. Learned representations of F0 and delexicalised speech have also proven useful for classifying persuasion and sarcasm [26]. These works demonstrate that prosodic correlates exhibit systematicity but leave open the question of at what time scale.

3. Experiments

Inspired by masked language models for text, we introduce a Masked Prosody Model (MPM) that learns to reconstruct corrupted sequences of pitch, loudness and voice activity during pretraining [12]. We investigate the effects of different corruption strategies on the utility of resulting representations across a set of downstream tasks that depend on prosodic features at varying timescales. While text units are much more intuitive to define than speech units, mask size plays an important role in the quality of text representations [27]. Similarly, we expect temporal and feature granularity to impact the complexity and value of MPM training. We use linear probes to compare MPM representations to hierarchical CWT encodings of acoustic correlates of prosody. Additionally, we use Conformer probes to compare the value of generic and task-dependent structures in both prosodic correlates and the full speech signal.

3.1. Masked Prosody Model

The architecture of the Masked Prosody Model, shown in Figure 1, consists of separate input-output sequences for each of the prosodic features and Conformer blocks [28] which are well-suited to their high resolution and continuous nature. Pitch (F0) and voice activity are extracted using the WORLD vocoder [29]. Energy is computed as the Root-Mean-Square (RMS) of each Mel Spectrogram frame. All features are extracted at a resolution of $\approx 10\text{ms}$ and normalised across the utterance to produce feature contours capturing important perceptual cues [30, 31]. Though normalisation removes information about e.g., the relative feature values for individual speakers, it allows encoding of unseen speakers. Before masking, input sequences are quantized into respective codebooks P , E and V of size c . Random segments of the aligned sequences, each with a mask length m , are masked until $(50 \pm 5)\%$ of the input signal remains. The model is trained to reconstruct each feature with an independent Categorical Cross Entropy loss. In training, feature losses are normalized by $1/\log c$ and their sum is optimized.²

²We release the code and pretrained weights for the MPM model at github.com/minixc/masked-prosody-model

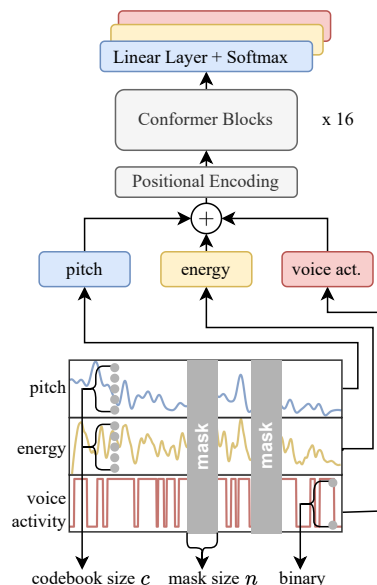


Figure 1: Architecture of the Masked Prosody Model.

To evaluate model representations, features from the 8th MPM Conformer layer are extracted as intermediary network representations are robust across tasks [32]. When aggregating at word or utterance levels, we compute the mean and maximum features for the target unit and concatenate them. Representations are evaluated using both linear probing [33] by training a linear classifier, and a more powerful Conformer classification model with 2 layers.

3.2. Downstream Tasks

We investigate prosodic structure by predicting human annotations known to be at least partially driven by prosodic features at different temporal scales. First, we assess whether representations capture local structure using the task of **syllable segmentation**. Syllabic units are often defined by changes to acoustic features including intensity and pitch [34] and are fundamental to models of speech perception; for example, infants without prior linguistic knowledge and adults in artificial language learning experiments are thought to rely on syllable-level units to acquire language [35, 36]. Second, prosody is known to support the linguistic organisation of speech above words through the demarcation of structure and salience [37, 16]. As such, we test **prominence and break detection** where CWT features have already proved useful. Finally, we investigate utterance-level systematicity using **emotion classification**, as affective information can be conveyed through prosody at the utterance level [38]. Rather than absolute performance, we are interested in the *relative* predictive power of different representations.

We evaluate **syllable segmentation** using **TIMIT**, a set of high-quality recordings from 630 speakers each reading ten phonetically-rich sentences [39]. For performing syllabification on **TIMIT**, previous works [40, 41] have framed the task as predicting vowel positions as a proxy for the number of syllables. Evaluation is done using Speaking Rate Error Rate (SER), which is defined as the absolute difference between actual and predicted syllables, divided by the actual number of syllables. The second metric is the correlation coefficient (Corr.) between the actual and predicted numbers of syllables per utterance.

Table 1: Probe evaluations using prosody representations (MPM, CWT, and untransformed representations of pitch, energy, and VAD) and speech representations (mel spectrogram, wav2vec, and HuBERT) as input. TIMIT syllable segmentation is evaluated with syllable error rate (SeR) and the correlation between true and predicted number of syllables; BURNc tasks with F1; RAUDESS emotion classification with weighted (WA) and unweighted accuracy (UA). Best task performance for prosody and speech models are in **bold**.

(a) Downstream task performance with the **linear probe**.

Type	Input Feature Name	Mask	TIMIT		BURNc (F1)		RAUDESS	
			SeR↓	Corr.↑	Boundary	Prominence	WA	UA
Speech	HuBERT	–	12.9	0.90	0.32	0.59	0.63	0.65
	Wav2Vec	–	12.8	0.85	0.28	0.56	0.54	0.55
	Mel Spectrogram	–	18.3	0.75	0.18	0.54	0.23	0.21
Prosody	Masked Prosody Model	4	14.9	0.85	0.27	0.51	0.15	0.16
		16	15.3	0.86	0.28	0.58	0.17	0.17
		128	16.0	0.78	0.25	0.50	0.21	0.22
		random	15.7	0.81	0.27	0.57	0.24	0.23
	Pitch, Energy, VAD (CWT)	–	19.8	0.71	0.02	0.39	0.20	0.19
Pitch, Energy, VAD	–	23.3	0.62	0.07	0.49	0.10	0.09	

(b) Downstream task performance with the **conformer probe**.

Type	Input Feature Name	TIMIT		BURNc (F1)		RAUDESS	
		SeR↓	Corr.↑	Boundary	Prominence	WA	UA
Speech	HuBERT	8.7	0.96	0.52	0.61	0.65	0.63
	Wav2Vec	10.2	0.95	0.51	0.60	0.63	0.63
	Mel Spectrogram	16.0	0.79	0.49	0.61	0.42	0.43
Prosody	Masked Prosody Model (random)	14.1	0.82	0.53	0.65	0.37	0.36
	Pitch, Energy, VAD (CWT)	16.1	0.78	0.46	0.51	0.31	0.29
	Pitch, Energy, VAD	17.0	0.76	0.44	0.50	0.22	0.21

ance. Given that syllables are marked by local rhythmic modulations in intensity and pitch, we don’t expect MPMs with wide masks to provide much benefit over untransformed input features [42, 43].

Prominence and break detection are evaluated on **BURNc** (the Lab News portion of Boston University Radio News Corpus) [44]. In this corpus, prominence is annotated with ToBI pitch accent types and boundaries between words are labelled for strength on a scale $\{0, \dots, 4\}$ [45]. We follow previous works and convert both tasks to binary classification where a word is prominent if any constituent syllable carries an accent ($\{H^*, L^*, L^*+H, L+H^*, H+, !H^*\}$) and boundaries consist of the strongest break indexes ($\{3,4\}$) [21].

There is no established train and test split for **BURNc**, and previous works use differing filtering and data selection methodologies [46, 21]. Due to this and the small dataset size, we use 5-fold cross validation and use the full dataset, including noisy samples. Using a handcrafted classification method over CWT features of pitch and energy, [21] report F1 respective performances of 0.85 and 0.59 for prominence and boundary detection. Pitch, energy, and durational information are useful, particularly in combination, for both tasks [47, 48]. We expect the additional flexibility of MPM encoding to be beneficial.

For **emotion classification**, we use the **RAUDESS** dataset (Ryerson Audio-Visual Database of Emotional Speech and Song Dataset), 1,440 utterances spoken by 24 speakers in eight emotions [49]. We select RAUDESS as emotion labels were validated with a wide battery of listening tests. Though we know of no purely prosodic baselines on this dataset, handcrafted prosodic features [50] and neural representations of delexicalised speech [26] have proved useful for emotion classification. This dataset consists of emotional renditions of two neutral phases. Though lexical content doesn’t affect our proposed rep-

resentations, this feature enables interesting comparisons with representations of the full speech signal.

3.3. Experimental Setup

We define MPM corruption strategies as combinations of mask m and codebook c sizes and expect the utility of self-supervision to be bounded by these parameters; e.g., small masks may be predicted trivially from continuous acoustic features, while recovering large masks will become impossible. In preliminary experiments, $m \in \{1, 2, 4, 8, 16, 32, 64, 128\}$ and $c \in \{4, 8, 16, 32, 64, 128, 512\}$ were tested; the smallest mask sizes consistently performed poorly and the best performance for all tasks was achieved with $c = 128$. For brevity, we report on strategies with $c = 128$ and a representative range $m \in \{4, 16, 128\}$. Motivated by [51] who find that masking randomly-sized spans of text yields more useful representations than masking individual tokens, we propose *random masking* as a generic corruption strategy; this involves sampling m uniformly between $(1, 128)$ for each batch. Using each corruption strategy, we train an MPM on the full LibriTTS dataset of over 500 hours of audiobook data [52]. Each model is trained for 10k steps (≈ 8 epochs, 30 minutes of TPUv3 runtime) with a batch size of 256, over utterances truncated to a maximum length of 6 seconds.

To enable comparison, CWT representations are generated from the same input features as MPMs (see Section 3.1) and no hyperparameter tuning or architecture adaptations are performed for downstream tasks. For topline results including the full speech signal, we use Wav2Vec Base [14], HuBERT Base [15] and Mel Spectrograms. For Wav2Vec and HuBERT, we use the model checkpoints trained on ≈ 900 h of LibriSpeech, which is comparable to the amount of data seen by the MPM.

The linear and Conformer (2 Conformer blocks, 5M pa-

rameters) classification probes are trained for 1000 steps with a batch size of 32, the AdamW optimizer, and a linear learning rate schedule with a maximum learning rate of $4e-5$ and 100 warmup steps.

4. Results

4.1. Effect of Corruption Strategy

In the MPM section of Table 1a, we use a linear probe to directly compare the performance achieved from models with $c = 128$ and $m \in \{4, 16, 128\}$ across the downstream tasks. We find that syllable segmentation benefits from smaller mask sizes; larger mask sizes may not encode sufficiently fine-grained local structures. Though boundary detection shows less sensitivity to mask size than prominence detection, the midsized mask performs best for both. Emotion recognition improves steadily with mask size, suggesting that some relevant features are only present when large mask sizes are encountered during training. The random-mask model performs well across all tasks. Although some local information useful for boundary detection and syllable segmentation may be degraded, it outperforms all corruption strategies for emotion recognition, indicating that random masking encodes both global and local structures and that emotion recognition benefits from this combination.

4.2. Comparing Structural Encodings of Prosody

We investigate the value of structure at different timescales by comparing MPM random-mask representations to the untransformed input features and CWT hierarchical representations.

Using **linear probes**, Table 1a shows the random-mask MPM consistently outperforms the untransformed and CWT features; the greatest relative gains of MPM are for BURNC boundary detection. The CWT prosodic features do not outperform the untransformed features for the BURNC tasks. This could indicate that some of the local features needed for both tasks are obscured by this transformation. Although both CWT and MPM outperform the untransformed input features for emotion recognition, MPM representations provide additional predictive power. This suggests that while hierarchical information is useful for this task, the additional structures learned through SSL are more valuable. MPM only slightly outperforms the input features for TIMIT syllable segmentation, confirming our expectation that this task doesn't require wide contextualization. For tasks that are expected to involve longer-ranging temporal context—prominence detection in BURNC and the RAVDESS emotion classification—the random mask model performs best overall.

In the Prosody section of Table 1b, we use more powerful **Conformer probes** to test whether the self-supervised MPM encodings provide value over learnable task-specific structures. Compared to the linear probe, the added capacity of the conformer produces higher performances across all representations; however, the magnitude of gains varies between tasks. TIMIT syllable segmentation detection shows relatively small improvements, indicating that the capacity of the probe was not a bottleneck on MPM performance. Boundary and prominence detection and emotion classification show larger improvements, suggesting that additional task-specific structures can be learnt by the Conformer probe. Though both CWT and MPM structures outperform the untransformed input for RAVDESS emotion classification, the flexibility of MPM is more valuable. The hierarchical structures encoded by CWT may not be sufficiently expressive for this more abstract task. Interestingly, though there are differences between their performance using the lin-

ear probe, the untransformed input features and CWT encodings achieve comparable performance for syllable segmentation and boundary and prominence detection with the Conformer probe. The hierarchical structures encoded by CWT only offer slight benefits over the task-specific structures learned from untransformed input features. MPM representations outperform both representations, again suggesting that useful information is more accessible in MPM features.

4.3. Investigating Additional Information in Speech

We also compare the SSL representations of the acoustic correlates of prosody to those of the full speech signal in Table 1b. Doing so highlights what information may be useful for particular downstream tasks. We might expect speech representations to consistently outperform MPM which only has access to a subset of information in the speech signal; however, Table 1b shows that MPM offers slight improvements over both wav2vec and HuBERT for BURNC boundary and prominence detection using the Conformer probe. This could indicate that salient information is more accessible in MPM representations. Although the SSL speech representations achieve the best TIMIT syllable segmentation performance, MPM features outperform the mel spectrogram: though phonetic information is valuable for such segmentation, this result highlights that syllable structure in TIMIT is conveyed through multiple speech features. All representations of the full speech signal outperform the prosody-only representations for emotion classification.

RAVDESS utterances are realisations of the same underlying lexical content [53], and therefore this cannot be explained by additional lexical content present in the full speech data. RAVDESS performance therefore depends on other features of speech beyond pitch, intensity, and voice activity, which could be explored in future work.

5. Discussion & Conclusions

We find that self-supervised methods can produce useful representations of the acoustic correlates of prosody, indicating that prosody exhibits predictable structure—systematicity— independently of lexical content. Our comparisons of corruption strategies reveal structure across timescales: syllable segmentation benefits from small masks while larger masks produce effective representations for detecting phrasal boundaries and emotion. By comparing MPM to hierarchical CWT representations of prosody, we find evidence for additional structural complexity; the flexibility of SSL is particularly valuable for emotion classification which involves complex dependencies across timescales. Our random-mask strategy encodes both long- and short-term structures, offering a generic representation of prosody that performs well across all tested tasks. To the best of our knowledge, we are the first to quantify the value of encoding structure of prosodic correlates isolated from lexical/segmental content at varying temporal scales using SSL mechanisms.

Given that the functions of prosody are tightly coupled with lexical content, it is unsurprising that representations of the full speech signal surpass MPM for tasks like emotion recognition. However, MPM is competitive for phrasal boundary and prominence detection, underscoring the value of prosody for these tasks. In the future, we hope to investigate the relative importance of pitch and energy features in this model, and more explicit inclusion of duration information. MPM also provides a means to compare prosodic systematicity across speech styles.

6. References

- [1] J. Hale, "A probabilistic Earley parser as a psycholinguistic model," in *NAACL*, 2001.
- [2] M. Brown, A. P. Salverda, L. C. Dilley, and M. K. Tanenhaus, "Metrical expectations from preceding prosody influence perception of lexical stress," *JEPHPP*, 2015.
- [3] S. Wallbridge, P. Bell, and C. Lai, "Quantifying the perceptual value of lexical and non-lexical channels in speech," in *Interspeech*, 2023.
- [4] M. Ip and A. Cutler, "Intonation facilitates prediction of focus even in the presence of lexical tones," in *Interspeech*, 2017.
- [5] F. Grosjean, "How long is the sentence? Prediction and prosody in the online processing of language," in *Linguistics*, 1983.
- [6] E. Ekstedt and G. Skantze, "How much does prosody help turn-taking? investigations using voice activity projection models," in *SIGDIAL*, 2022.
- [7] A. Cutler and J. McQueen, "How prosody is both mandatory and optional," in *Above and beyond the segments*, 2014.
- [8] C. Lai, "What do you mean, you're uncertain? The interpretation of cue words and rising intonation in dialogue," in *Interspeech*, 2010.
- [9] J. Mills, "Delexicalised auditory priming of implicit prosody," in *Speech Prosody*, 2020.
- [10] A. Cervone, C. Lai, S. Pareti, and P. Bell, "Towards automatic detection of reported speech in dialogue using prosodic cues," in *Interspeech*, 2015.
- [11] J. Donahue *et al.*, "DeCAF: A deep convolutional activation feature for generic visual recognition," in *PMLR*, 2014.
- [12] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *NAACL*, 2019.
- [13] D. Ofer, N. Brandes, and M. Linial, "The language of proteins: NLP, machine learning & protein sequences," *CSBJ*, 2021.
- [14] S. Schneider, A. Baevski, R. Collobert, and M. Auli, "wav2vec: Unsupervised pre-training for speech recognition," *arXiv:1904.05862*, 2019.
- [15] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, "HuBERT: Self-supervised speech representation learning by masked prediction of hidden units," *TASLP*, 2021.
- [16] M. de Seyssel *et al.*, "ProsAudit, a prosodic benchmark for self-supervised speech models," in *Interspeech*, 2023.
- [17] R. Skerry-Ryan *et al.*, "Towards end-to-end prosody transfer for expressive speech synthesis with Tacotron," in *ICML*, 2018.
- [18] L. Qu, T. Li, C. Weber, T. Pékarek-Rosin, F. Ren, and S. Wermter, "Disentangling prosody representations with unsupervised speech reconstruction," *TASLP*, 2023.
- [19] G. Ioannides, M. Owen, A. Fletcher, V. Rozgic, and C. Wang, "Towards paralinguistic-only speech representations for end-to-end speech emotion recognition," in *Interspeech*, 2023.
- [20] A. T. Sigurgeirsson and S. King, "Do prosody transfer models transfer prosody?" in *ICASSP*, 2023.
- [21] A. Suni, J. Simko, D. Aalto, and M. Vainio, "Hierarchical representation and estimation of prosody using continuous wavelet transform," *Comput. Speech Lang.*, 2017.
- [22] A. Suni, S. Kakouros, M. Vainio, and J. Simko, "Prosodic prominence and boundaries in sequence-to-sequence speech synthesis," in *Speech Prosody*, 2020.
- [23] E. O. Selkirk, *Phonology and Syntax: The Relation between Sound and Structure*. MIT Press, 1984.
- [24] M. Beckman, *Stress And Non-Stress Accent*. De Gruyter, 1986.
- [25] S. Kakouros and O. Räsänen, "Automatic detection of sentence prominence in speech using predictability of word-level acoustic features," in *Interspeech*, 2015.
- [26] J. Weston, R. Lenain, U. Meepegama, and E. Fristed, "Learning de-identified representations of prosody from raw audio," in *ICML*, 2021.
- [27] Y. Levine *et al.*, "PMI-Masking: Principled masking of correlated spans," in *ICLR*, 2020.
- [28] A. Gulati *et al.*, "Conformer: Convolution-augmented transformer for speech recognition," in *Interspeech*, 2020.
- [29] M. Morise, F. Yokomori, and K. Ozawa, "WORLD: a vocoder-based high-quality speech synthesis system for real-time applications," *IEICE Transactions*, 2016.
- [30] A. Rosenberg, "Classification of prosodic events using quantized contour modeling," in *NAACL*, 2010.
- [31] A. E. Kimball and J. Cole, "Pitch contour shape matters in memory," in *ICSP*, 2016.
- [32] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson, "How transferable are features in deep neural networks?" in *NeurIPS*, 2014.
- [33] G. Alain and Y. Bengio, "Understanding intermediate layers using linear classifier probes," *ICLR Workshop*, 2017.
- [34] S. G. Parker, *Quantifying the sonority hierarchy*. University of Massachusetts Amherst, 2002.
- [35] P. W. Jusczyk, L. J. Kennedy, and A. M. Jusczyk, "Young infants' retention of information about syllables," *Infant Behavior and Development*, 1995.
- [36] T. Matzinger, N. Ritt, and W. T. Fitch, "The influence of different prosodic cues on word segmentation," *Front. in Psychology*, 2021.
- [37] J. Cole, "Prosody in context: A review," *Language, Cognition and Neuroscience*, 2015.
- [38] A. Batliner, S. Steidl, D. Seppi, and B. Schuller, "Segmenting into adequate units for automatic recognition of emotion-related episodes: a speech-based approach," *AHCI*, 2010.
- [39] J. S. Garofolo, "TIMIT acoustic phonetic continuous speech corpus," *LDS*, 1993.
- [40] Y. Jiao, V. Berisha, M. Tu, and J. Liss, "Convex weighting criteria for speaking rate estimation," *TASLP*, 2015.
- [41] J. Yuan, N. Ryant, X. Cai, K. Church, and M. Liberman, "Automatic recognition of suprasegmentals in speech," *arXiv:2108.01122*, 2021.
- [42] O. Räsänen, G. Doyle, and M. C. Frank, "Pre-linguistic segmentation of speech into syllable-like units," *Cognition*, vol. 171, 2018.
- [43] O. Bagou, C. Fougerson, and U. Frauenfelder, "Contribution of prosody to the segmentation and storage of words" in the acquisition of a new mini-language," in *Speech Prosody*, 2002.
- [44] M. Ostendorf, P. J. Price, and S. Shattuck-Hufnagel, "The Boston University radio news corpus," *LDS*, 1995.
- [45] K. Silverman *et al.*, "ToBI: a standard for labeling English prosody," in *SLP*, 1992.
- [46] S. Ananthakrishnan and S. S. Narayanan, "Automatic prosodic event detection using acoustic, lexical, and syntactic evidence," *TASLP*, 2007.
- [47] B. Ludusan and E. Dupoux, "Towards low-resource prosodic boundary detection," in *SLTU*, 2014.
- [48] O. Kalinli and S. S. Narayanan, "A saliency-based auditory attention model with applications to unsupervised prominent syllable detection in speech," in *Interspeech*, 2007.
- [49] S. R. Livingstone and F. A. Russo, "The Ryerson audio-visual database of emotional speech and song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English," *PLOS ONE*, 2018.
- [50] H. Cao, Š. Beňuš, R. C. Gur, R. Verma, and A. Nenkova, "Prosodic cues for emotion: analysis with discrete characterization of intonation," in *Speech Prosody*, 2014.
- [51] M. Joshi, D. Chen, Y. Liu, D. S. Weld, L. Zettlemoyer, and O. Levy, "SpanBERT: Improving pre-training by representing and predicting spans," *TACL*, 2020.
- [52] H. Zen *et al.*, "LibriTTS: A Corpus Derived from LibriSpeech for Text-to-Speech," in *Interspeech*, 2019.
- [53] L. Tian, J. Moore, and C. Lai, "Recognizing emotions in spoken dialogue with hierarchically fused acoustic and lexical features," in *SLT*, 2016.