

Muffin: Muffled Audio Encoding with Filter-Based Masking

Marcel A. Vélez Vásquez¹

¹Music Cognition Group, ILLC, University of Amsterdam

Abstract

Masked autoencoders have advanced representation learning in vision and language, yet audio remains dominated by spectrogram-based approaches which then are treated like images, disregarding audio-specific characteristics. We propose *Muffled Audio Encoding*, a framework for self-supervised learning directly on raw waveforms using 1D transformers and masking through time-domain filters (e.g., low-, high-, and band-pass). This approach encourages representations that capture long-range and frequency-selective dependencies without requiring Fourier transforms and losing phase information. We outline our design and experimental plan for evaluating this method across multiple audio domains.

1 Introduction

In recent years, self-supervised pretraining tasks have become central to learning effective representations for language [1], images [2], and audio [3]. Among them, *masked prediction* tasks such as BERT’s masked language modeling [4] and Masked Image Modeling (MIM) [5] have shown that reconstructing missing information from context can yield transferable features.

However, masked modeling has not yet been fully adapted to the unique structure of audio. Existing waveform-based methods, such as Wav2Vec 2.0 [6], typically mask contiguous segments in the *time domain*, which helps the model learn temporal dependencies but ignores rich frequency relationships that define many audio phenomena. Spectrogram-based approaches [7, 8] instead operate in the time–frequency plane but often treat spectrograms as if they were ordinary images. However, unlike images, where both spatial axes represent the same modality, spectrograms have fundamentally different axes: time and frequency. Consequently, such models may overlook the asymmetric structure of audio data and, in most cases, disregard phase information that is crucial for faithful signal reconstruction. Audio signals exhibit complex dependencies across both time and frequency. By masking in the *frequency domain*, models can learn to reconstruct missing spectral information and capture phenomena such as harmonics, timbral structure, and spectral envelopes—crucial for music, speech, and environmen-

tal sounds. To exploit this, we propose a method that performs masking directly in the waveform domain through audio **filters** (low-, high-, band-pass, and band-stop), effectively suppressing selected frequency regions while maintaining phase continuity.

We term this approach *MUFFIN* (**M**uffled audio **e**ncoding), which bridges time- and frequency-domain masking through filter-based perturbations within a masked autoencoding framework. We hypothesize that frequency-selective masking fosters more general and perceptually grounded representations, transferable across tasks such as speech recognition, music analysis, and environmental sound classification.

2 Related work

Self-supervised learning has reshaped audio representation learning by leveraging masked or contrastive objectives without explicit labels. Most approaches adapt ideas from language or vision domains but differ in what they mask and how they represent audio.

Masked prediction in language and vision. BERT [4] introduced masked language modeling with low masking ratios (15%) for discrete tokens. Vision Masked Autoencoders (MAE) [5] extended this idea to continuous signals, masking up to 75% of image patches and reconstructing missing pixels with an asymmetric encoder–decoder design. The MAE formulation emphasizes spatial redundancy and efficient pretraining.

Time-domain masking for speech. **Wav2Vec 2.0** [6] learns contextualized speech embeddings by predicting quantized latent representations of masked time segments. The model is trained on raw waveforms but relies on contrastive objectives and masking entire sequences. **WavLM** [9] extends this to multi-task pretraining with denoising.

Spectrogram-domain masking. **Mockingjay** [7] apply masked prediction on spectrograms, typically masking 10–20% of time–frequency patches. These methods benefit from the structured representation of spectrograms but discard phase information and introduce fixed spectral resolution.

091 The more recent **Masked Spectrogram Mod-** **Evaluation and transfer learning.** Evaluation 138
 092 **eling (MSM-MAE)** [8] adapts vision MAE to is conducted through **transfer learning** on the 139
 093 2D spectrogram inputs, showing strong transfer to **Jukemir**[12] benchmark, which provides standard- 140
 094 HEAR 2021 tasks. ized tasks across musical attributes: 141

095 **Bridging domains.** While waveform-based mod-
 096 els retain temporal precision, they typically mask
 097 only contiguous time spans. Spectrogram-based
 098 models can mask across both time and frequency but
 099 often disregard phase information and treat both
 100 axes as the same modality. Our method bridges
 101 these perspectives by *masking directly in the wave-*
 102 *form domain through filtering*: rather than masking
 103 discrete spectrogram patches, we apply parametric
 104 low-, high-, band-pass, or band-stop filters as struc-
 105 tured masks. This removes specific frequency bands
 106 in a continuous manner while preserving tempo-
 107 ral coherence and encouraging representations that
 108 capture both long-range and frequency-selective de-
 109 pendencies.

110 3 method

111 We employ a similar architecture to the Vision trans-
 112 former, where 2D convolution-layers are replaced by
 113 1D ones. Each input waveform segment is randomly
 114 *muffled* using filters drawn from, or a combination
 115 of:

- 116 • low-pass: remove high-frequency content;
- 117 • high-pass: remove low-frequency content;
- 118 • or, band-pass / band-stop: isolate or suppress
 119 mid-range content

120 The model is trained to reconstruct the original
 121 waveform from these masked variants, promoting
 122 robustness to spectral variation and encouraging
 123 multi-scale temporal representations.

124 4 Experimental Setup

125 **Pretraining and study design.** We begin with
 126 pretraining on the **MagnaTagATune (MTT)**[10]
 127 dataset to systematically explore the effect of **mask-**
 128 **ing ratio, decoder depth and width, and filter**
 129 **sampling strategies.** These ablations test whether
 130 filter-based masking can yield competitive represen-
 131 tations on modest data scales. We include basic
 132 augmentations such as random gain, polarity in-
 133 version, and time-stretching to evaluate robustness.
 134 After validating configurations on MTT, we scale
 135 pretraining to **AudioSet**[11](2M 10-second clips) to
 136 study how representation quality scales with dataset
 137 size and domain diversity.

- **MTT**[10]: music tagging (multi-label classifi- 142
 cation) 143
- **GiantSteps**[13]: musical key estimation 144
- **GTZAN**[14]: Genre classification 145
- **EmoMusic**[15]: valence–arousal regression 146

We further include **OpenMIC**[16] & **NSynth**[17] 147
 for instrument recognition. 148

Cross-domain evaluation. Beyond Jukemir[12], 149
 we plan comparisons on broader audio domains: 150

- **Environmental sounds:** ESC-50[18] and Ur- 151
 banSound8K (US8K)[19] 152
- **Speech:** SPCV2[20], VC1[21], and CREMA- 153
 D[22] for speaker and emotion recognition 154

Testing whether frequency-selective masking gener- 155
 alizes beyond music. 156

Baselines and comparisons. We compare 157
 against a range of pretrained self-supervised and 158
 contrastive models: 159

- **Waveform-based:** Wav2Vec 2.0[6], 160
 WavLM[9], ATST-Base[23], TUNe[24] 161
- **Spectrogram-based:** MSM-MAE[8], 162
 M2D[25], MATPAC++[26] 163
- **Retrieval-based:** SLAP[27] and Golden Re- 164
 triever[28] 165

All baselines are only compared on the datasets of 166
 the original papers. 167

5 Future Work 168

We will next carry out the experiments outlined 169
 above, implementing MUFFIN with large-scale pre- 170
 training and executing the full set of comparisons 171
 across self-supervised, contrastive, and retrieval- 172
 based baselines. These experiments will analyze the 173
 influence of filter type, masking ratio, and model 174
 capacity on representation quality and examine 175
 whether these factors behave differently for audio 176
 than they do in vision-based masked autoencoding. 177
 We will also evaluate MUFFIN representations on 178
 the HEAR challenge to assess general-purpose audio 179
 performance, but due to the two-page limit, we omit 180
 the full list of evaluation tasks here. 181

182 **References**

- 183 [1] Y. Gu, R. Tinn, H. Cheng, M. Lucas, N.
184 Usuyama, X. Liu, T. Naumann, J. Gao, and
185 H. Poon. “Domain-specific language model
186 pretraining for biomedical natural language
187 processing”. In: *ACM Transactions on Com-
188 puting for Healthcare (HEALTH)* 3.1 (2021),
189 pp. 1–23.
- 190 [2] T. Chen, S. Kornblith, M. Norouzi, and G.
191 Hinton. “A simple framework for contrastive
192 learning of visual representations”. In: *Interna-
193 tional conference on machine learning*. PMLR.
194 2020, pp. 1597–1607.
- 195 [3] J. Spijkervet and J. A. Burgoyne. “Contrastive
196 learning of musical representations”. In: *arXiv
197 preprint arXiv:2103.09410* (2021).
- 198 [4] J. Devlin, M.-W. Chang, K. Lee, and K.
199 Toutanova. “Bert: Pre-training of deep bidirec-
200 tional transformers for language understand-
201 ing”. In: *Proceedings of the 2019 conference of
202 the North American chapter of the association
203 for computational linguistics: human language
204 technologies, volume 1 (long and short papers)*.
205 2019, pp. 4171–4186.
- 206 [5] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and
207 R. Girshick. “Masked autoencoders are scal-
208 able vision learners”. In: *Proceedings of the
209 IEEE/CVF conference on computer vision and
210 pattern recognition*. 2022, pp. 16000–16009.
- 211 [6] A. Baevski, Y. Zhou, A. Mohamed, and M.
212 Auli. “wav2vec 2.0: A framework for self-
213 supervised learning of speech representations”.
214 In: *Advances in neural information processing
215 systems* 33 (2020), pp. 12449–12460.
- 216 [7] A. T. Liu, S.-w. Yang, P.-H. Chi, P.-c. Hsu, and
217 H.-y. Lee. “Mockingjay: Unsupervised speech
218 representation learning with deep bidirectional
219 transformer encoders”. In: *ICASSP 2020-2020
220 IEEE International Conference on Acoustics,
221 Speech and Signal Processing (ICASSP)*. IEEE.
222 2020, pp. 6419–6423.
- 223 [8] D. Niizumi, D. Takeuchi, Y. Ohishi, N. Harada,
224 and K. Kashino. “Masked spectrogram mod-
225 eling using masked autoencoders for learn-
226 ing general-purpose audio representation”. In:
227 *HEAR: Holistic Evaluation of Audio Repre-
228 sentations*. PMLR. 2022, pp. 1–24.
- 229 [9] S. Chen, C. Wang, Z. Chen, Y. Wu, S. Liu, Z.
230 Chen, J. Li, N. Kanda, T. Yoshioka, X. Xiao,
231 et al. “Wavlm: Large-scale self-supervised pre-
232 training for full stack speech processing”. In:
233 *IEEE Journal of Selected Topics in Signal Pro-
234 cessing* 16.6 (2022), pp. 1505–1518.
- [10] E. Law, K. West, M. Mandel, M. Bay, and
J. Downie. “Evaluation of algorithms using
games: the case of music annotation”. In: *Pro-
ceedings of the 11th International Society for
Music Information Retrieval Conference (IS-
MIR)*. Utrecht, the Netherlands. 2010.
- [11] J. F. Gemmeke, D. P. Ellis, D. Freedman, A.
Jansen, W. Lawrence, R. C. Moore, M. Plakal,
and M. Ritter. “Audio set: An ontology and
human-labeled dataset for audio events”. In:
*2017 IEEE international conference on acous-
tics, speech and signal processing (ICASSP)*.
IEEE. 2017, pp. 776–780.
- [12] R. Castellon, C. Donahue, and P. Liang. “Codi-
fied audio language modeling learns useful rep-
resentations for music information retrieval”.
In: (2021).
- [13] P. Knees, Á. Faraldo Pérez, P. Herrera Boyer,
R. Vogl, S. Böck, F. Hörschläger, and M. Le
Goff. “Two data sets for tempo estimation
and key detection in electronic dance music
annotated from user corrections”. In: *ISMIR*
(2015).
- [14] G. Tzanetakis and P. Cook. “Musical genre
classification of audio signals”. In: *IEEE
Transactions on speech and audio processing*
10.5 (2002), pp. 293–302.
- [15] M. Soleymani, M. N. Caro, E. M. Schmidt,
C.-Y. Sha, and Y.-H. Yang. “1000 songs for
emotional analysis of music”. In: *Proceedings
of the 2nd ACM international workshop on
Crowdsourcing for multimedia*. 2013, pp. 1–6.
- [16] E. Humphrey, S. Durand, and B. McFee.
“OpenMIC-2018: An Open Data-set for Multi-
ple Instrument Recognition.” In: *ISMIR*. 2018,
pp. 438–444.
- [17] J. Engel, C. Resnick, A. Roberts, S. Dieleman,
M. Norouzi, D. Eck, and K. Simonyan. “Neural
audio synthesis of musical notes with wavenet
autoencoders”. In: *International conference
on machine learning*. PMLR. 2017, pp. 1068–
1077.
- [18] K. J. Piczak. “ESC: Dataset for environmen-
tal sound classification”. In: *Proceedings of the
23rd ACM international conference on Multi-
media*. 2015, pp. 1015–1018.
- [19] J. Salamon, C. Jacoby, and J. P. Bello. “A
dataset and taxonomy for urban sound re-
search”. In: *Proceedings of the 22nd ACM in-
ternational conference on Multimedia*. 2014,
pp. 1041–1044.
- [20] P. Warden. “Speech Commands: A Dataset
for Limited-Vocabulary Speech Recognition”.
In: *CoRR* (2018).

- 289 [21] A. Nagrani, J. S. Chung, and A. Zisserman.
290 “VoxCeleb: A Large-Scale Speaker Identifica-
291 tion Dataset”. In: *Proc. Interspeech 2017*. 2017,
292 pp. 2616–2620.
- 293 [22] H. Cao, D. G. Cooper, M. K. Keutmann, R. C.
294 Gur, A. Nenkova, and R. Verma. “Crema-d:
295 Crowd-sourced emotional multimodal actors
296 dataset”. In: *IEEE transactions on affective
297 computing* 5.4 (2014), pp. 377–390.
- 298 [23] X. Li, N. Shao, and X. Li. “Self-supervised au-
299 dio teacher-student transformer for both clip-
300 level and frame-level tasks”. In: *IEEE/ACM
301 Transactions on Audio, Speech, and Language
302 Processing* 32 (2024), pp. 1336–1351.
- 303 [24] M. A. V. Vásquez, J. A. Burgoyne, et al.
304 “Tailed U-Net: Multi-Scale Music Representa-
305 tion Learning.” In: *ISMIR*. 2022, pp. 67–75.
- 306 [25] D. Niizumi, D. Takeuchi, Y. Ohishi, N. Harada,
307 and K. Kashino. “Masked modeling duo:
308 Learning representations by encouraging both
309 networks to model the input”. In: *ICASSP
310 2023-2023 IEEE International Conference
311 on Acoustics, Speech and Signal Processing
312 (ICASSP)*. IEEE. 2023, pp. 1–5.
- 313 [26] A. Queleñec, P. Chouteau, G. Peeters, and
314 S. Essid. “MATPAC++: Enhanced Masked
315 Latent Prediction for Self-Supervised Audio
316 Representation Learning”. In: *arXiv preprint
317 arXiv:2508.12709* (2025). URL: [https://
318 arxiv.org/abs/2508.12709](https://arxiv.org/abs/2508.12709).
- 319 [27] J. Guinot, A. Riou, E. Quinton, and G.
320 Fazekas. “SLAP: Siamese Language-Audio
321 Pretraining Without Negative Samples for Mu-
322 sic Understanding”. In: (2025).
- 323 [28] J. Guinot, E. Quinton, and G. Fazekas. “GD-
324 Retriever: Controllable Generative Text-Music
325 Retrieval with Diffusion Models”. In: (2025).