

# SUPERVISED AND UNSUPERVISED LEARNING OF AUDIO REPRESENTATIONS FOR MUSIC UNDERSTANDING

M. C. McCallum, F. Korzeniowski,  
S. Oramas, F. Gouyon, A. F. Ehmann

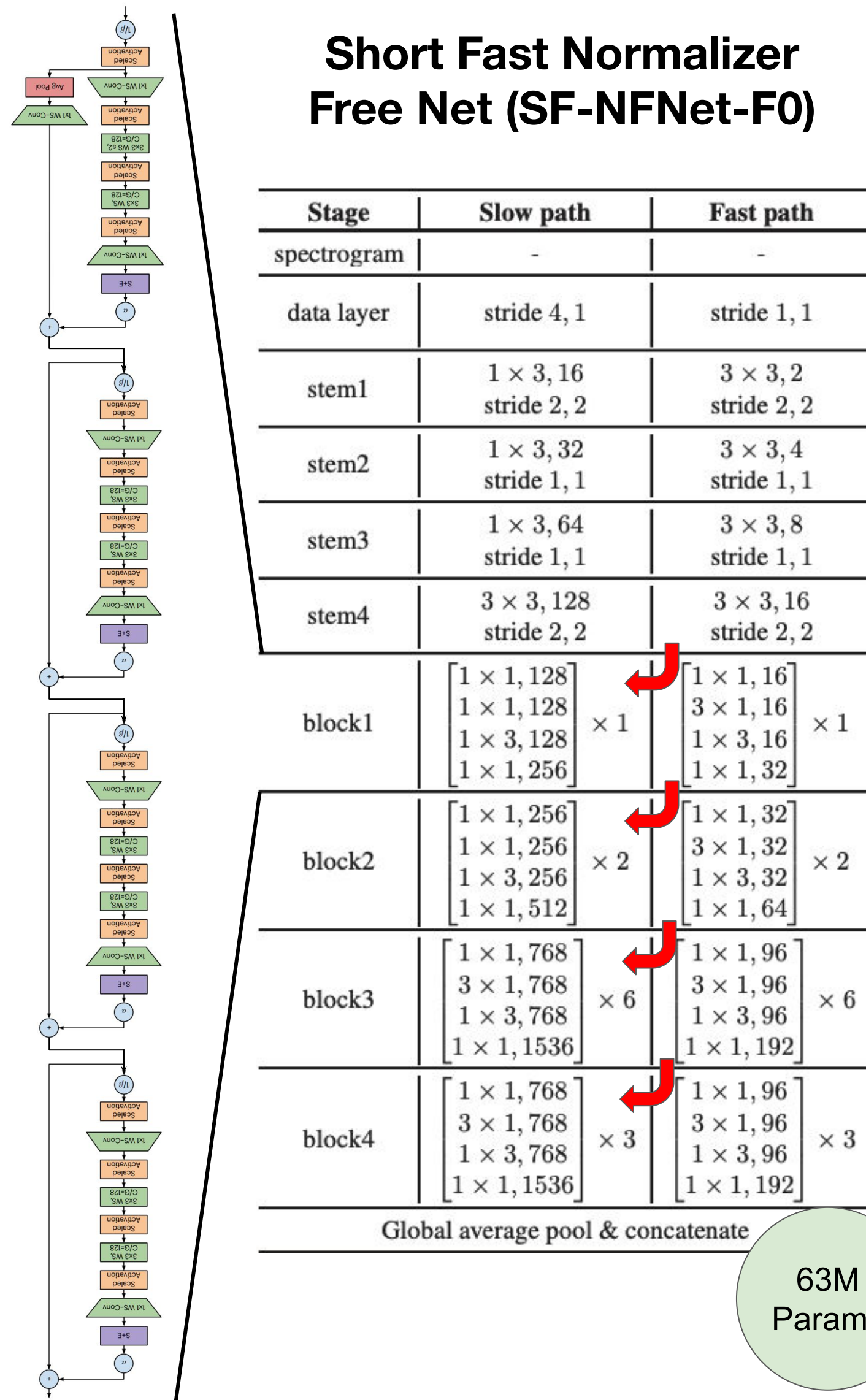
Sirius XM, USA



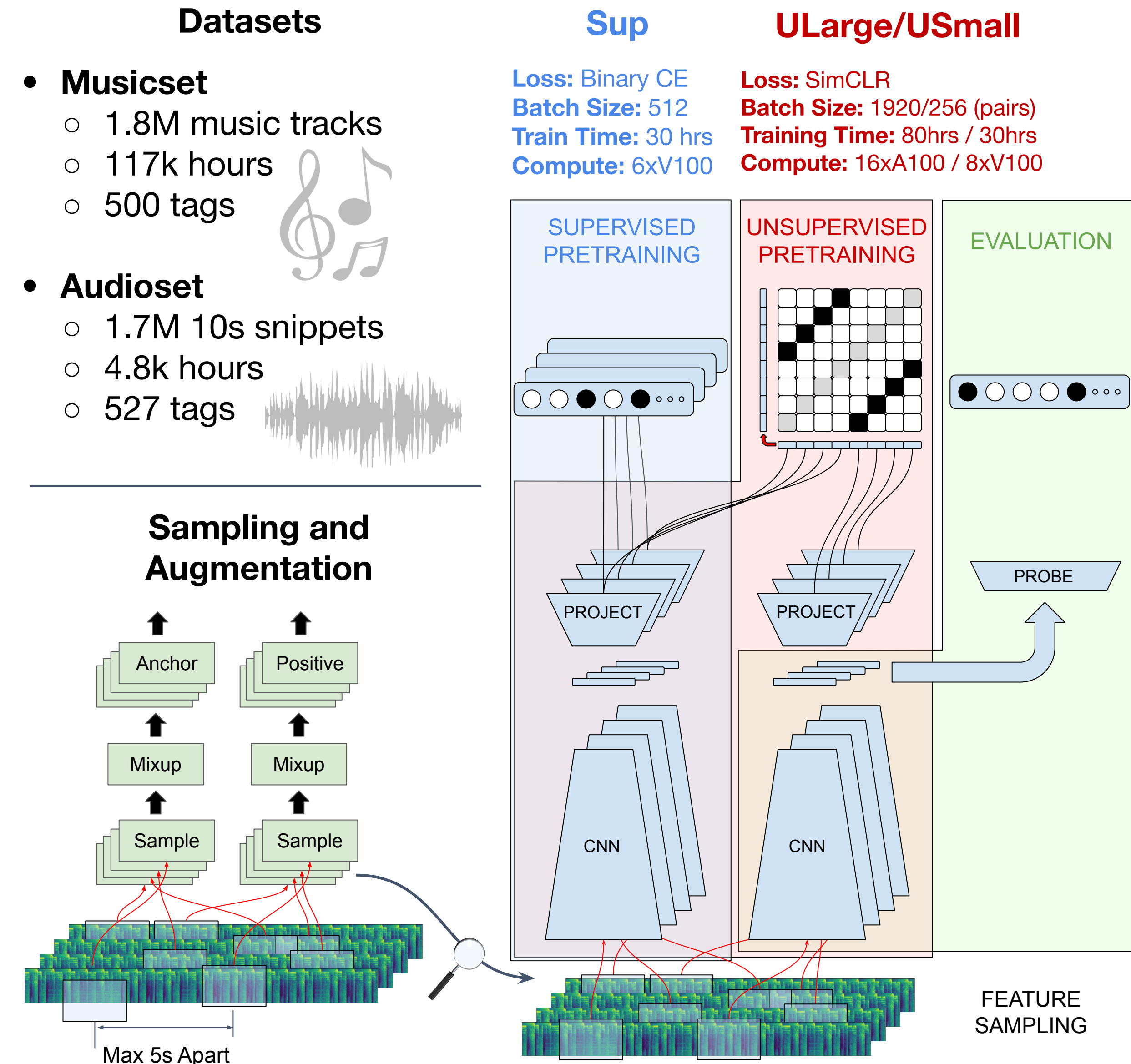
## OBJECTIVES

- Provide a **broad set of baselines** for music understanding tasks
- **Compare** the effectiveness of **supervised and unsupervised learning objectives at scale**
- Investigate the **impact of training dataset content** and batch size for training unsupervised models
- Release a model to **enable and accelerate downstream research** in audio and / or multimodal understanding for music.

## MODEL

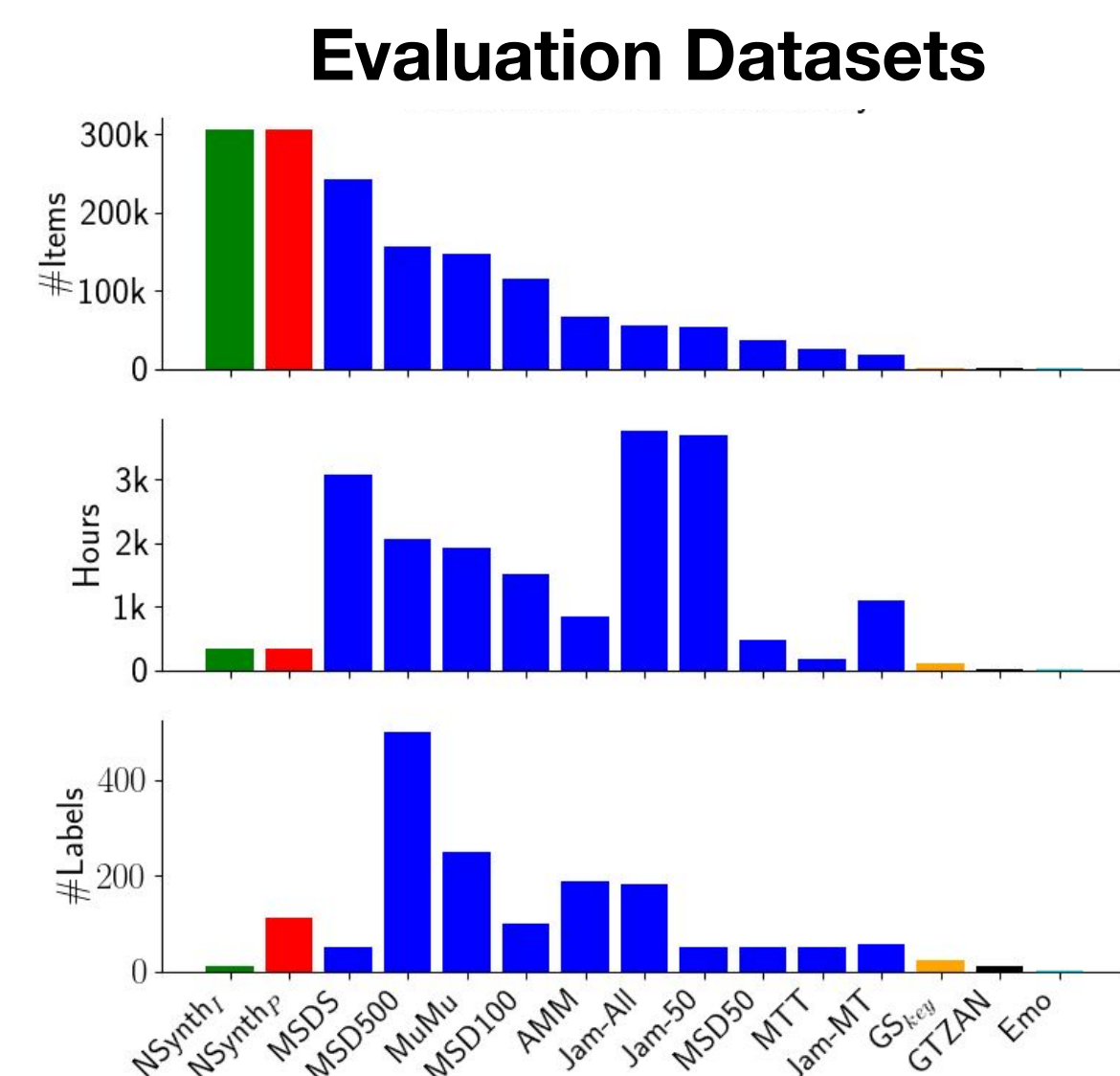


## PRETRAINING METHODOLOGY



## EVALUATION

- **7 Distinct Audio Collections**
- **15 Datasets / Annotations**
- Embeddings **global-average pooled** along track length.
- Probes consist of **MLPs**
- Probe **hyperparameters optimized**, respecting **same restrictions** as previous audio representation work



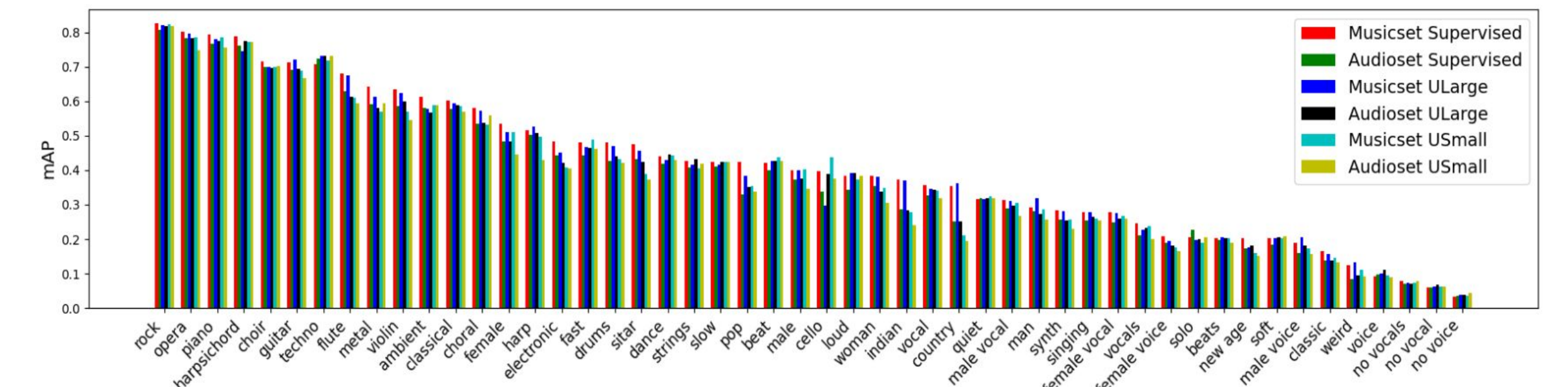
## RESULTS

Model	MSDS		MSD50		MSD100		MSD500		MuMu		AMM		Jam-MT	
	mAP	ROC	mAP	ROC	mAP	ROC	mAP	ROC	mAP	ROC	mAP	ROC	mAP	ROC
Musicset-Sup	<b>0.363</b>	<b>0.903</b>	<b>0.459</b>	<b>0.913</b>	<b>0.346</b>	<b>0.906</b>	<b>0.169</b>	<b>0.898</b>	<b>0.257</b>	<b>0.908</b>	<b>0.180</b>	<b>0.791</b>	<b>0.161</b>	<b>0.786</b>
Audioset-Sup	0.308	0.880	0.375	0.883	0.278	0.877	0.128	0.874	0.191	0.867	0.156	0.760	0.137	0.749
Musicset-ULarge	0.351	<b>0.900</b>	0.438	<b>0.908</b>	0.321	<b>0.897</b>	0.152	<b>0.891</b>	0.235	0.893	<b>0.174</b>	<b>0.784</b>	<b>0.158</b>	<b>0.781</b>
Audioset-ULarge	0.311	0.885	0.377	0.886	0.276	0.878	0.121	0.873	0.162	0.855	0.156	0.763	0.142	0.765
Musicset-USmall	0.319	0.888	0.384	0.892	0.283	0.881	0.129	0.878	0.190	0.871	0.155	0.762	0.138	0.757
Audioset-USmall	0.286	0.876	0.353	0.878	0.251	0.870	0.110	0.868	0.152	0.850	0.151	0.753	0.136	0.753
SOTA	0.348	<b>0.897</b>	0.386	<b>0.921</b>	0.185	-	-	-	-	0.888*	0.163	0.773	<b>0.161</b> †	<b>0.781</b> †
	[15]	[15]	[14]	[14]	[22]	-	-	-	-	[42]	[37]	[37]	[49]	[49]

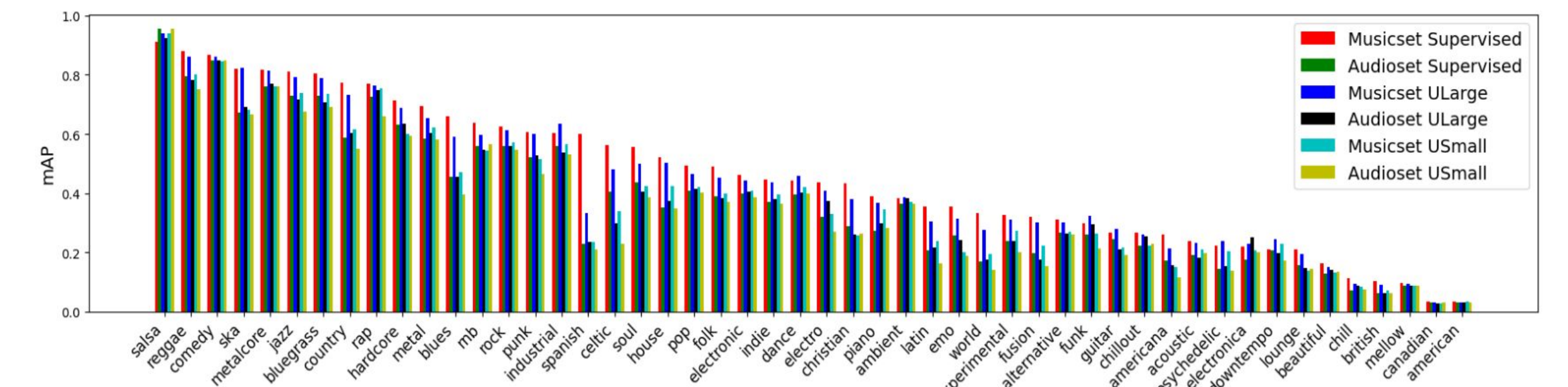
  

Model	MTT		GTZAN	NSynth	NSynth <sub>1</sub>	Emo <sub>v</sub>	Emo <sub>A</sub>	GS <sub>Key</sub>	Jam-50		Jam-All	
	mAP	ROC	Acc	Acc	Acc	r <sup>2</sup>	r <sup>2</sup>	W. Acc	mAP	ROC	mAP	ROC
Musicset-Sup	<b>0.413</b>	<b>0.917</b>	<b>0.835</b>	0.793	0.731	0.566	<b>0.726</b>	0.286	<b>0.321</b>	<b>0.843</b>	<b>0.162</b>	<b>0.839</b>
Audioset-Sup	0.386	0.904	0.748	0.819	0.676	0.341	0.545	0.210	0.284	0.822	0.135	0.813
Musicset-ULarge	<b>0.404</b>	<b>0.914</b>	0.735	<b>0.892</b>	0.740	0.577	0.700	0.667	<b>0.317</b>	<b>0.839</b>	<b>0.159</b>	<b>0.833</b>
Audioset-ULarge	0.391	0.906	0.672	0.805	0.721	0.438	0.624	0.287	0.285	0.826	0.131	0.816
Musicset-USmall	0.389	0.905	0.686	0.824	0.714	0.389	0.668	0.508	0.292	0.828	0.138	0.817
Audioset-USmall	0.375	0.897	0.648	0.777	0.698	0.386	0.609	0.197	0.268	0.817	0.127	0.809
Jukebox [23]	<b>0.414</b>	<b>0.915</b>	0.797	-	-	<b>0.617</b>	<b>0.721</b>	0.667	-	-	-	-
Prev. SF-NFNet-F0 [2]	0.395	-	-	0.880	<b>0.782</b>	-	-	-	-	-	-	-
SOTA Excl. [2, 23]	0.384	<b>0.92</b>	0.821	-	0.741	0.556	0.704	<b>0.796*</b>	0.298	0.832	-	-
	[37]	[12]	[11]	-	[43]	[44]	[45]	[46]	[47]	[47]	-	-

## Magnatagatune



## MSD50



## KEY TAKE-AWAYS

- Supervised models achieve SotA on all multilabel tagging tasks
- Unsupervised models generalize better to novel tasks like pitch and key
- Music understanding models **perform better** when pretrained on purely music data

**Musicset-ULarge Model Available Here:**  
<https://github.com/PandoraMedia/music-audio-representations>

