

Self-Similarity-Based and Novelty-Based Loss for Music Structure Analysis

Geoffroy Peeters LTCl, Télécom-Paris, IP-Paris
https://github.com/geoffroypeeters/ssmnet_ISMIR2023

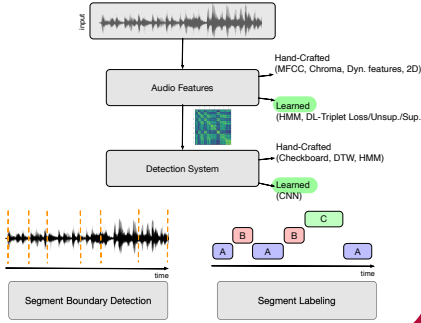
Introduction

Music structure analysis ?

- identify musical segments that compose a music track (also known as Segment Boundary estimation)
- label them based on their similarity (also known as segment labelling).

Over the years and the accessibility of annotated datasets, systems have switched from

- hand-crafted audio features used as input to hand-crafted detection systems
- hand crafted audio features used as input to deep learning detection
- deep learned features used as input to hand-crafted detection systems



Proposal

We propose a system which relies on deep learned audio features used as input to deep learning detection system

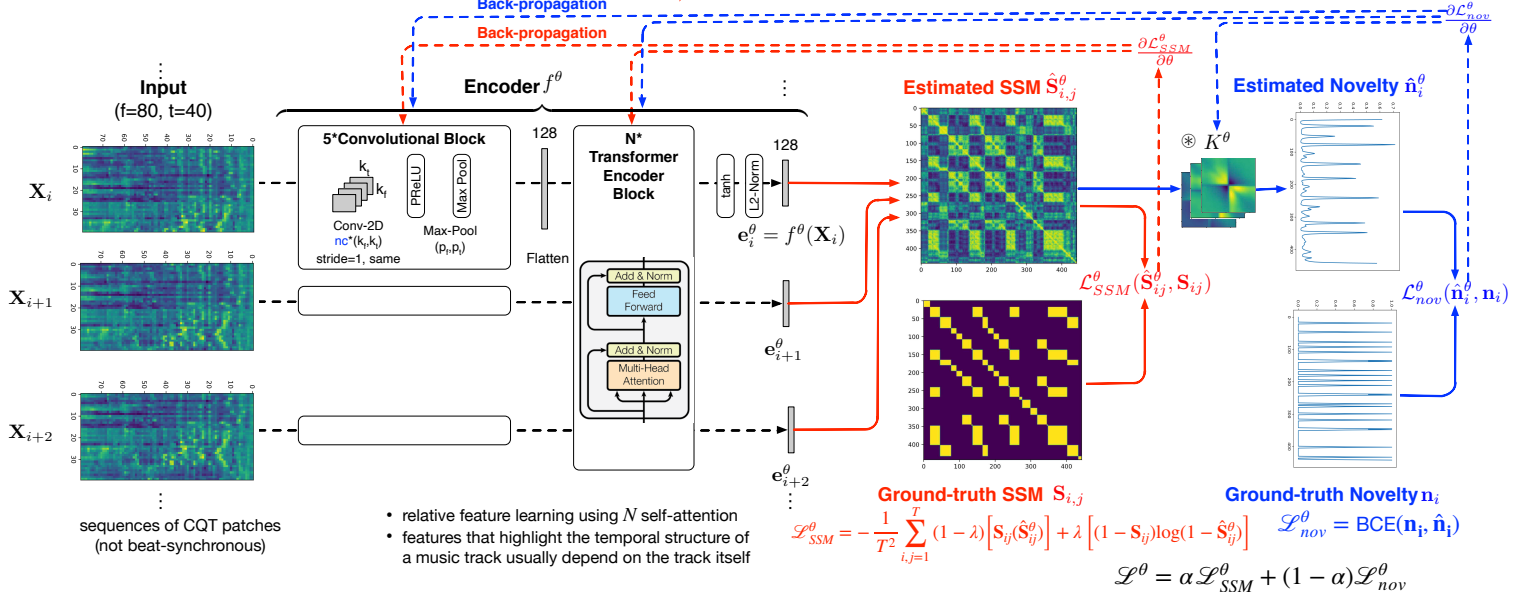
We propose a straightforward training paradigms based on:

- a direct comparison of estimated and ground-truth SSM
- a direct comparison of estimated and ground-truth novelty-curve
- We train a deep encoder f^θ and a set of convolution kernels K^θ such that
- the Self-Similarity-Matrix (SSM) \hat{S}_{ij}^θ resulting from the learned features $\mathbf{e}_i^\theta = f^\theta(\mathbf{X}_i)$ approximates a ground-truth SSM \mathbf{S}_{ij}
- the Novelty-curve $\hat{\mathbf{n}}_i$ resulting from convolving \hat{S}_{ij}^θ with learned kernels K^θ approximates a ground-truth Novelty \mathbf{n}_i

$$\frac{\partial \mathcal{L}_{SSM}^\theta}{\partial \theta} = \sum_{i,j=1}^T \frac{\partial \mathcal{L}_{SSM}^\theta}{\partial \hat{S}_{ij}^\theta} \left(\frac{\partial \hat{S}_{ij}^\theta}{\partial \mathbf{e}_i^\theta} \frac{\partial \mathbf{e}_i^\theta}{\partial \theta} + \frac{\partial \hat{S}_{ij}^\theta}{\partial \mathbf{e}_j^\theta} \frac{\partial \mathbf{e}_j^\theta}{\partial \theta} \right)$$

$$\hat{S}_{ij}^\theta = 1 - \frac{1}{4} \|\mathbf{e}_i^\theta - \mathbf{e}_j^\theta\|_2^2 \in [0,1]$$

- K^θ
- 3 kernels of (41,41) corresponds to 20s + (1x1) convolution
- Initialised using Foote ck kernels / using randn

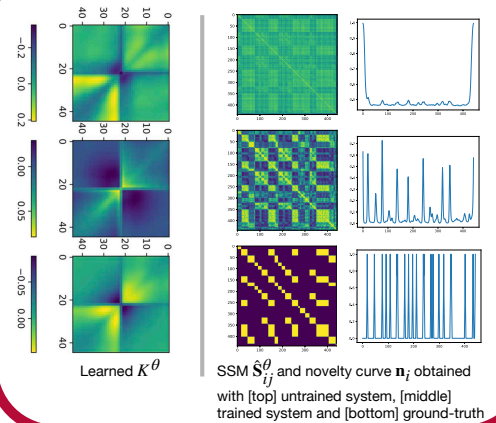


Evaluation

Previous results
 Our results
 Ablation study

	RWC-Pop-AIST		SA-Pop		SA-1A		SA-Two		Annotation
	HR.5F	HR.3F	HR.5F	HR.3F	HR.5F	HR.3F	HR.5F	HR.3F	
Grill [23, 42] G&I	.506	.715	-	-	-	-	.541	.623	Up./An.*
McCallum [25] Unsynch.	-	-	-	-	.497	.535	-	-	Up./An.*
Beat-synch.	-	-	-	-	-	-	.337	.563	Up./An.*
Salamon [30] DEFT-unsup/	.438	.653	.447	.623	-	-	.356	.553	Up./An.-1/2
Wang [27] cluster/D	-	.681	-	-	-	-	.597 / .595	.611 / .600	Low/An-1/2
Buisson [29] HE0/HE1	-	-	-	-	-	-	-	-	Low/An-1/2
Ours (best conf.)	.399	.713	.298 / .295	.631 / .624	.250 / .261	.520 / .511	.231 / .237	.521 / .530	Up./An-1/2
			.296 / .318	.570 / .610	.302 / .336	.547 / .612	.287 / .287	.589 / .589	
Ablation study N									
N=3/alpha=0.5/K:train-Init:chk		.713		.532		.472		.448	Up./An-1
N=2/alpha=0.5/K:train-Init:chk		.701		.535		.474		.449	Up./An-1
N=1/alpha=0.5/K:train-Init:chk		.677		.631		.520		.521	Up./An-1
N=0/alpha=0.5/K:train-Init:chk		.696		.535		.459		.443	Up./An-1
Ablation study alpha									
N=3/alpha=1/K:train-Init:chk	.154		.121		.102		.111		Up./An-1
N=3/alpha=0/K:train-Init:chk	.007		.120		.026		.095		Up./An-1
Ablation study K^theta									
N=3/alpha=0.5/K:train-Init:randn		.713		.543		.470		.457	Up./An-1
N=1/alpha=0.5/K:train-Init:randn		.709		.547		.470		.457	Up./An-1
N=3/alpha=0.5/K:fix-Init:chk		.330		.250		.199		.196	Up./An-1

Illustration



[23] T. Grill and J. Schlüter, "Music Boundary Detection Using Neural Networks on Combined Features and Two-Level Annotations," in Proc. of ISMIR, Malaga, Spain, 2015.
 [25] M. C. McCallum, "Unsupervised Learning of Deep Features for Music Segmentation," in Proc. of IEEE ICASSP Brighton, UK, May 2019.
 [27] T. Grill and J. Schlüter, "Structural segmentation with convolutional neural networks MIREX submission," 2015.
 [30] J. Salamon, O. Nieto, and N. J. Bryan, "Deep embeddings and section fusion improve music segmentation," in Proc. of ISMIR, Online, November, 8-12 2021.
 [29] J.-C. Wang, J. B. L. Smith, W.-T. Lu, and X. Song, "Supervised metric learning for music structure features," in Proc. of ISMIR, Online, 2021.
 [29] M. Buisson, B. McFee, S. Essid, and H.-C. Crayencour, "Learning multi-level representations for hierarchical music structure analysis," in Proc. of ISMIR, 2022.