

# Are deep neural networks really learning relevant features?

Corey Kereliuk

DTU Compute

Technical University of Denmark

cmke@dtu.dk

Bob L. Sturm

Audio Analysis Lab

Aalborg University

bst@create.aau.dk

Jan Larsen

DTU Compute

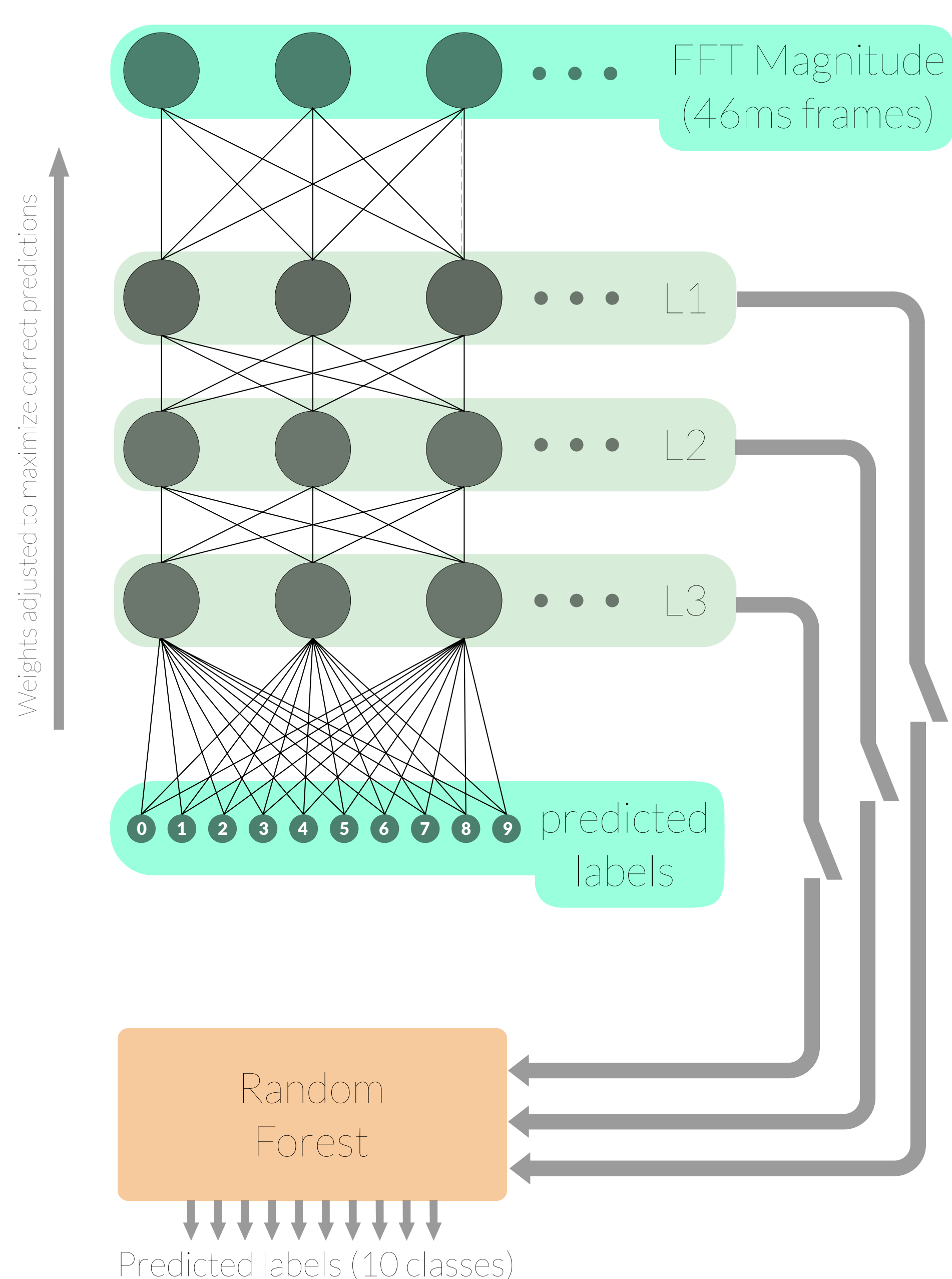
Technical University of Denmark

janla@dtu.dk

## ABSTRACT

Deep neural networks (DNNs) are being evaluated more and more for machine learning. This is due to a variety of factors, including improvements to training algorithm efficiency, the capacity of DNNs to implicitly learn features, and their excellent performance in many domains. Two recent works [1,2] apply DNNs to content-based music information retrieval, specifically music genre recognition (MGR). The conclusions in these works are unfortunately drawn from evaluations using the GTZAN dataset, which is now known to contain faults (replicated observations and artists) that have major effects when not taken into consideration [3]. We thus re-examine the conclusions in these works considering these faults, and are led to question the degree to which the learned features are actually an improvement over standard "hand-crafted" features such as Mel-frequency cepstral coefficients (MFCCs). All our results are reproducible with the open software repository: <https://github.com/coreyker/dnn-mgr>

## SETUP



## DNN architecture

- 46 ms FFT magnitude frames (513-d) input
- 50 or 500 rectified linear units (ReLU) per hidden layer
- 3 hidden layers
- Trained using SGD with (and without) 'dropout'

## Random forest

- Input L1, L2, L3, or all three layers concatenated
- Trained with (and without) temporal aggregation of frames
- 500 trees in forest

Table of GTZAN classification accuracy from [1]:

Hidden Units	Layer	No Aggregation		Aggregation	
		ReLU+SGD	ReLU+SGD+Dropout	ReLU+SGD	ReLU+SGD+Dropout
50	1	75.0 ± 1.3	75.0 ± 1.4	75.0 ± 1.7	76.5 ± 1.5
	2	75.4 ± 1.2	77.5 ± 2.2	79.6 ± 2.7	77.0 ± 2.2
	3	78.3 ± 1.1	77.0 ± 1.2	81.3 ± 1.8	78.0 ± 1.0
	All	79.0 ± 2.0	78.0 ± 1.6	81.3 ± 1.9	81.5 ± 1.7
500	1	72.7 ± 2.8	73.5 ± 1.9	71.8 ± 0.7	75.5 ± 1.1
	2	78.5 ± 1.9	78.5 ± 2.9	79.5 ± 1.9	82.5 ± 1.8
	3	80.5 ± 1.4	79.5 ± 2.6	83.0 ± 1.2	82.0 ± 1.4
	All	79.0 ± 1.4	80.5 ± 1.8	82.5 ± 2.3	83.0 ± 1.1

Results from our replication of systems and evaluation in [1]:

Hidden Units	Layer	No Aggregation		Aggregation	
		ReLU+SGD	ReLU+SGD+Dropout	ReLU+SGD	ReLU+SGD+Dropout
50	1	73.20	73.60	74.40	77.60
	2	74.40	73.60	79.20	77.20
	3	76.40	72.80	78.40	76.00
	All	75.20	75.60	81.60	76.80
500	1	67.20	72.80	71.60	76.40
	2	69.60	79.60	74.40	78.80
	3	72.00	81.60	76.40	80.00
	All	71.60	80.40	76.40	82.00

Results with consideration of faults in GTZAN:

Hidden Units	Layer	No Aggregation		Aggregation	
		ReLU+SGD	ReLU+SGD+Dropout	ReLU+SGD	ReLU+SGD+Dropout
50	1	40.34	40.00	39.66	39.66
	2	38.62	39.66	42.07	41.03
	3	38.28	40.00	44.48	38.62
	All	39.31	39.66	42.41	39.31
500	1	42.07	42.07	39.66	40.69
	2	41.72	43.79	41.03	45.52
	3	40.69	44.14	42.76	48.97
	All	40.69	43.79	41.03	47.59

Confusions of two systems (highlighted above) built and tested without (left) and with (right) consideration of GTZAN faults:

	blues	classical	country	disco	hiphop	jazz	metal	pop	reggae	rock	Pr
blues	96.0	0.0	8.0	4.0	0.0	0.0	4.0	0.0	8.0	4.0	77.4
classical	0.0	88.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	100.0
country	0.0	0.0	84.0	0.0	0.0	0.0	0.0	0.0	4.0	12.0	84.0
disco	0.0	0.0	0.0	72.0	4.0	0.0	0.0	8.0	0.0	8.0	78.3
hiphop	0.0	0.0	0.0	8.0	76.0	0.0	0.0	16.0	0.0	0.0	76.0
jazz	0.0	8.0	0.0	0.0	0.0	100.0	0.0	4.0	0.0	0.0	89.3
metal	0.0	0.0	0.0	0.0	16.0	0.0	88.0	0.0	4.0	0.0	81.5
pop	0.0	0.0	0.0	12.0	0.0	0.0	88.0	0.0	8.0	0.0	81.5
reggae	0.0	0.0	4.0	4.0	0.0	0.0	0.0	60.0	4.0	0.0	83.3
rock	4.0	4.0	4.0	0.0	4.0	0.0	8.0	4.0	4.0	64.0	66.7
F	85.7	93.6	84.0	75.0	76.0	94.3	84.6	84.6	69.8	65.3	81.6

	blues	classical	country	disco	hiphop	jazz	metal	pop	reggae	rock	Pr
blues	12.9	0.0	20.0	3.4	0.0	0.0	3.7	0.0	0.0	6.2	28.6
classical	0.0	100.0	0.0	0.0	3.7	18.5	0.0	0.0	0.0	0.0	83.8
country	12.9	0.0	16.7	3.4	0.0	0.0	0.0	0.0	0.0	12.5	35.7
disco	9.7	0.0	26.7	24.1	0.0	22.2	3.7	10.0	11.5	43.8	15.6
hiphop	3.2	0.0	0.0	17.2	55.6	0.0	3.7	13.3	26.9	3.1	44.1
jazz	32.3	0.0	13.3	0.0	3.7	33.3	3.7	6.7	7.7	3.1	30.0
metal	12.9	0.0	0.0	3.4	0.0	0.0	81.5	0.0	0.0	6.2	75.9
pop	0.0	0.0	3.3	3.4	25.9	14.8	0.0	63.3	15.4	0.0	52.8
reggae	3.2	0.0	6.7	24.1	11.1	0.0	0.0	3.3	34.6	18.8	31.0
rock	12.9	0.0	13.3	20.7	0.0	11.1	3.7	3.3	3.8	6.2	9.1
F	17.8	91.2	22.7	18.9	49.2	31.6	78.6	57.6	32.7	7.4	42.4

## Conclusion:

Our evaluation of the DNN systems in [1, 2] produces dramatically lower figures of merit when considering the faults in GTZAN. In other work [3], we find Bayesian classification using standard 'hand-crafted' features (MFCCs) produces a classification accuracy of about 50% in GTZAN when considering its faults. A random forest with MFCCs produces a classification accuracy of 51%. These results thus contradict the conclusion that these DNNs are producing features better or more relevant than "hand-crafted" MFCCs for MGR. We are led to the question, are DNNs really learning relevant features?

- [1] S. Sigita and S. Dixon, "Improved music feature learning with deep neural networks," ICASSP, 2014.
- [2] P. Hamel and D. Eck, "Learning features from music audio with deep belief networks." ISMIR, 2010.
- [3] B. L. Sturm, "The state of the art ten years after a state of the art: Future research in music information retrieval," J. New Music Research, 43(2), pp. 147--172, 2014.