

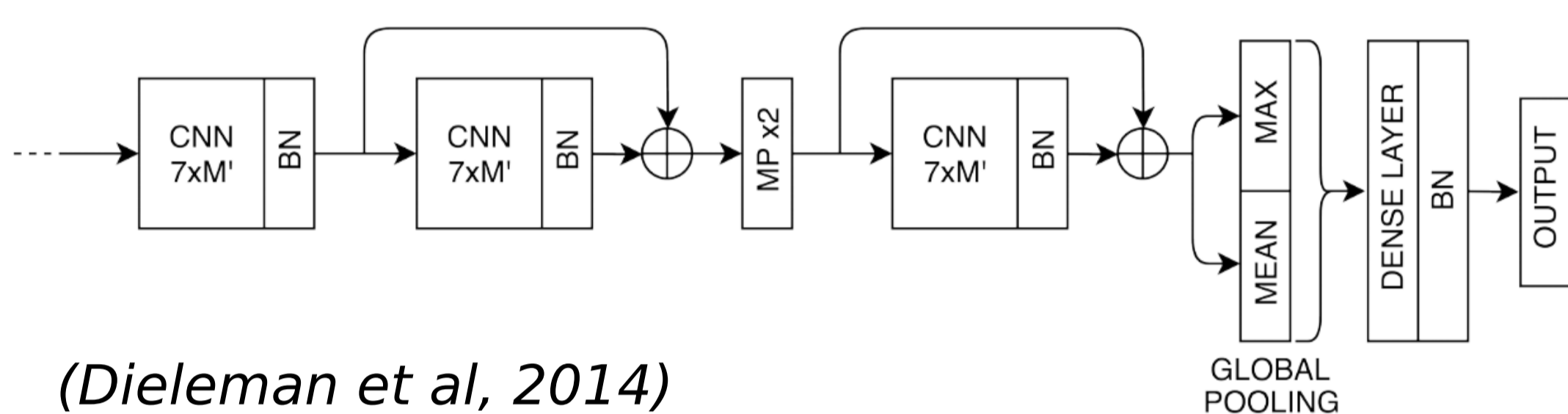
# End-to-end learning for music audio tagging at scale

## waveform > spectrogram?

can waveform-based deep learning models achieve better performance than spectrogram-based ones?

**YES!** If enough training data is available

### Shared back-end for a fair front-end comparison



(Dieleman et al, 2014)

### MagnaTagATune dataset: 25k songs

	ROC AUC	PR AUC	# param
<i>Waveform models</i>			
Dieleman et al.	85.58	29.59	194k
SampleCNN	88.56	34.38	2.4M
Waveform model (ours)	<b>89.05</b>	<b>34.92</b>	11.8M
Waveform model (smaller)	88.94	34.47	1.3M
<i>Spectrogram models</i>			
Timbre CNN	89.07	34.92	220k
VGG	89.99	37.56	450k
Spectrogram model (ours)	<b>90.40</b>	<b>38.11</b>	5M
Spectrogram model (smaller)	90.28	37.55	222k

few training data?

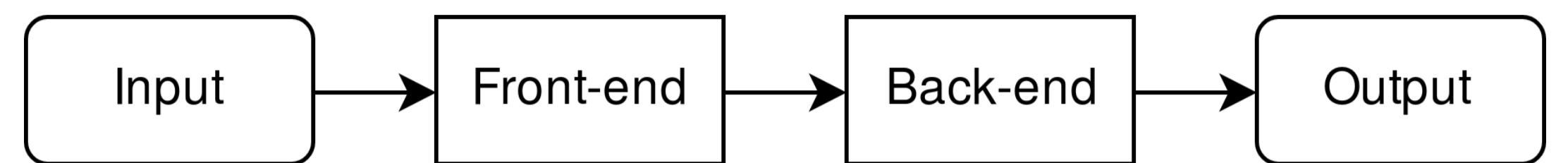
**spectrogram > waveform**

### Million Song Dataset: 250k songs

	ROC AUC	PR AUC	# param
<i>Waveform models</i>			
SampleCNN	88.12	-	2.4M
SampleCNN multi-level & multi-scale	<b>88.42</b>	-	-
Waveform model (ours)	87.41	28.53	5.3M
<i>Spectrogram models</i>			
VGG + RNN	86.2	-	3M
Multi-level & multi-scale	<b>88.78</b>	-	-
Spectrogram model (ours)	<b>88.75</b>	<b>31.24</b>	5.9M

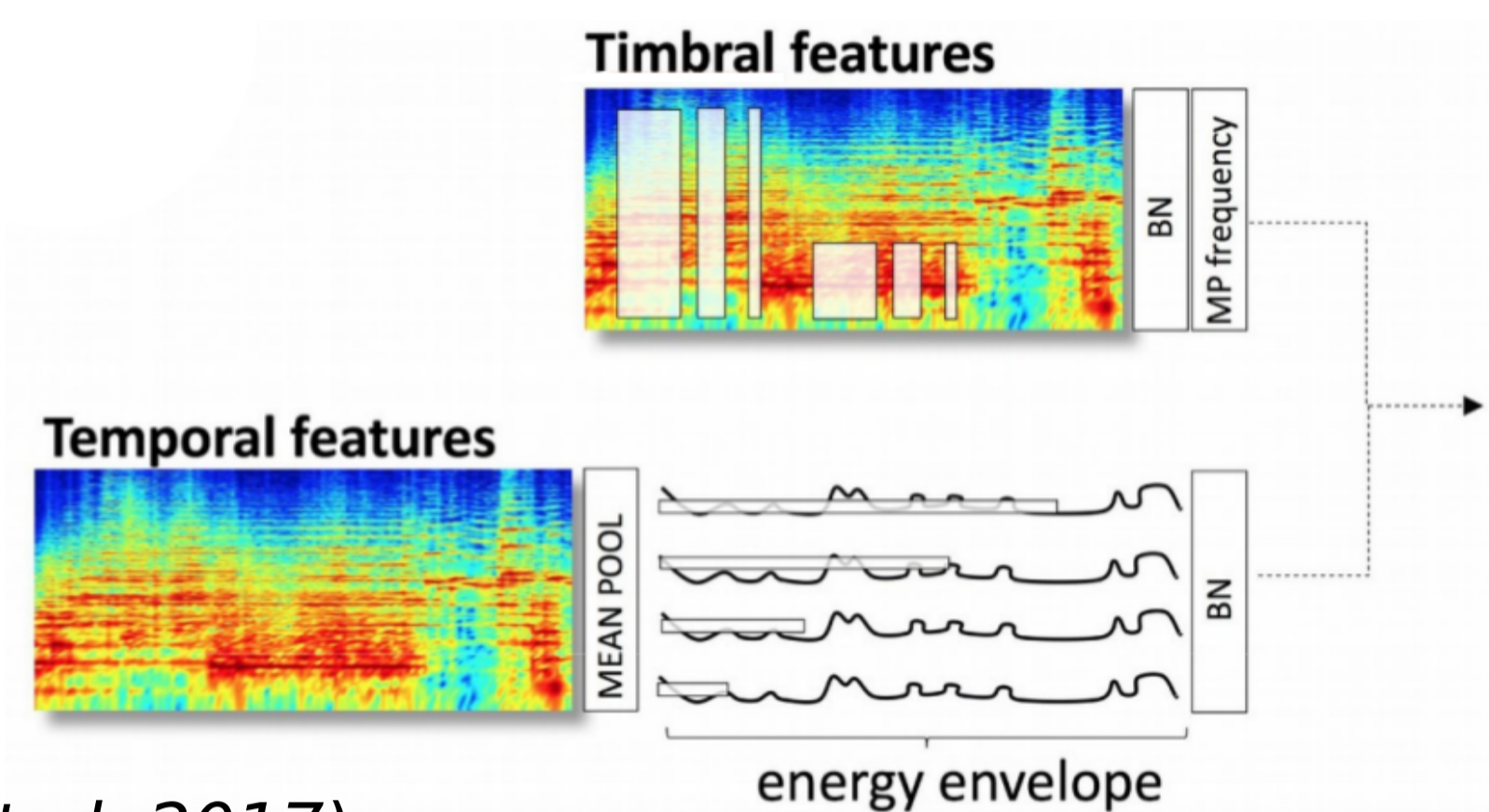
a reasonable amount of training data?

**spectrogram > waveform**



### Spectrogram front-end

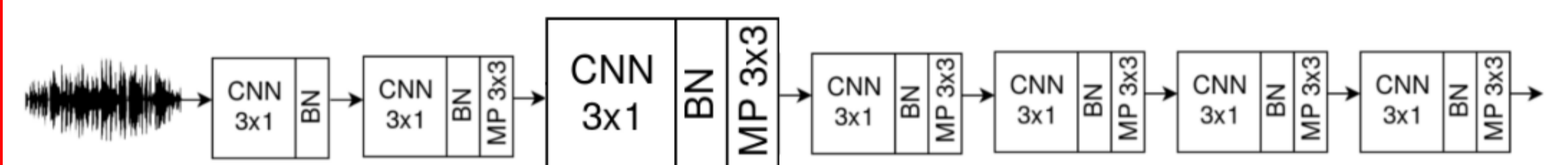
heavily based on domain knowledge



(Pons et al, 2017)

### Waveform front-end

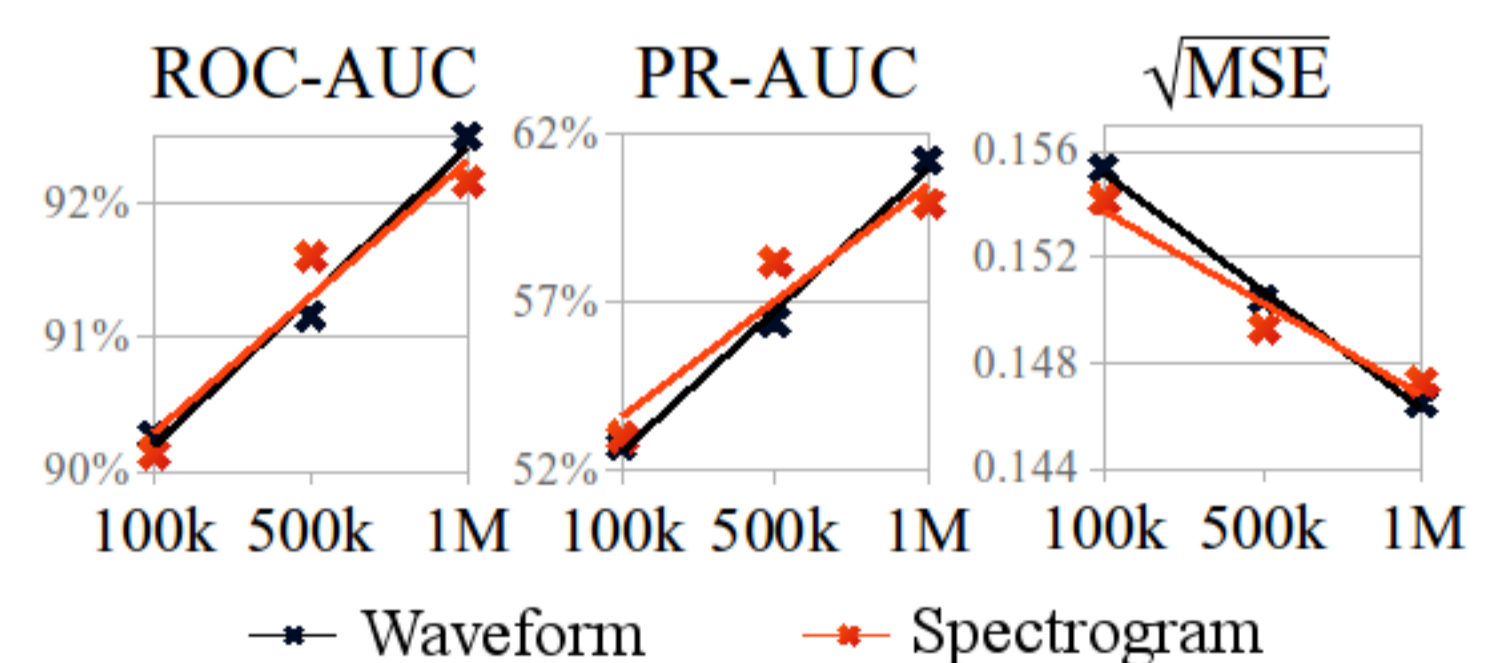
an assumption-free model



(Lee et al, 2017)

### Private dataset: 1.2M songs

Models	train size	ROC AUC	PR AUC	$\sqrt{MSE}$
Baseline	1.2M	91.61%	54.27%	0.1569
Waveform	1M	<b>92.50%</b>	<b>61.20%</b>	<b>0.1465</b>
Spectrogram	1M	92.17%	59.92%	0.1473
Waveform	500k	91.16%	56.42%	0.1504
Spectrogram	500k	91.61%	58.18%	0.1493
Waveform	100k	90.27%	52.76%	0.1554
Spectrogram	100k	90.14%	52.67%	0.1542



lots of training data?

**waveform > spectrogram**

### Qualitative results

Bias towards predicting popular tags  
"lead vocals", "English" or "male vocals"

Predicting each tag independently vs.  
predicting all tags together

"East Coast" vs. "West Coast"

"Baroque period" vs. "Classic period"