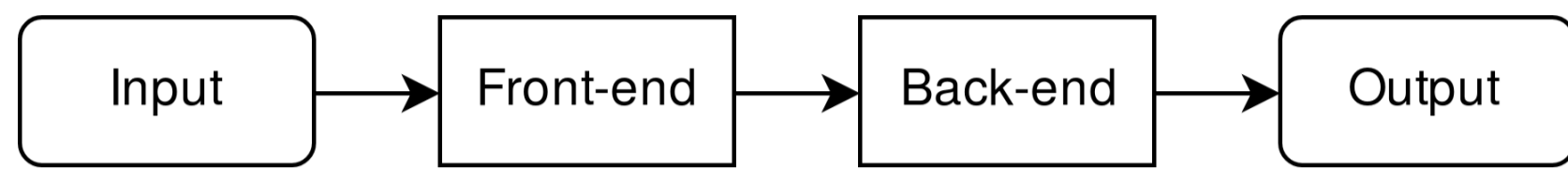


# End-to-end learning for music audio tagging at scale

Which deep learning architecture shall we use for audio tagging?

How much data is available?  
A human-annotated corpus of **1.5M songs**

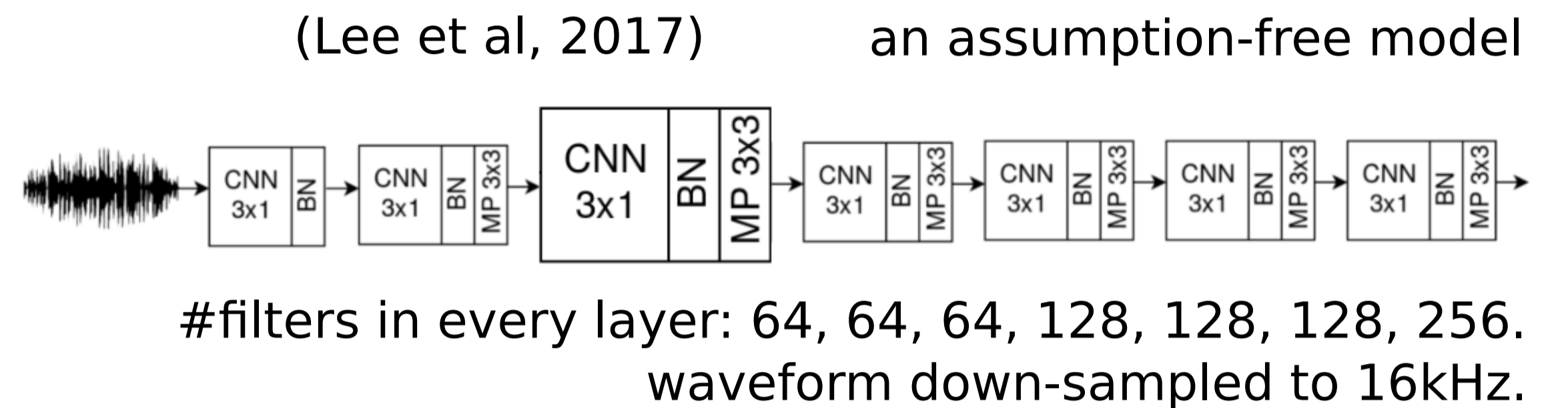


Can waveform front-ends achieve better performance than spectrogram front-ends?

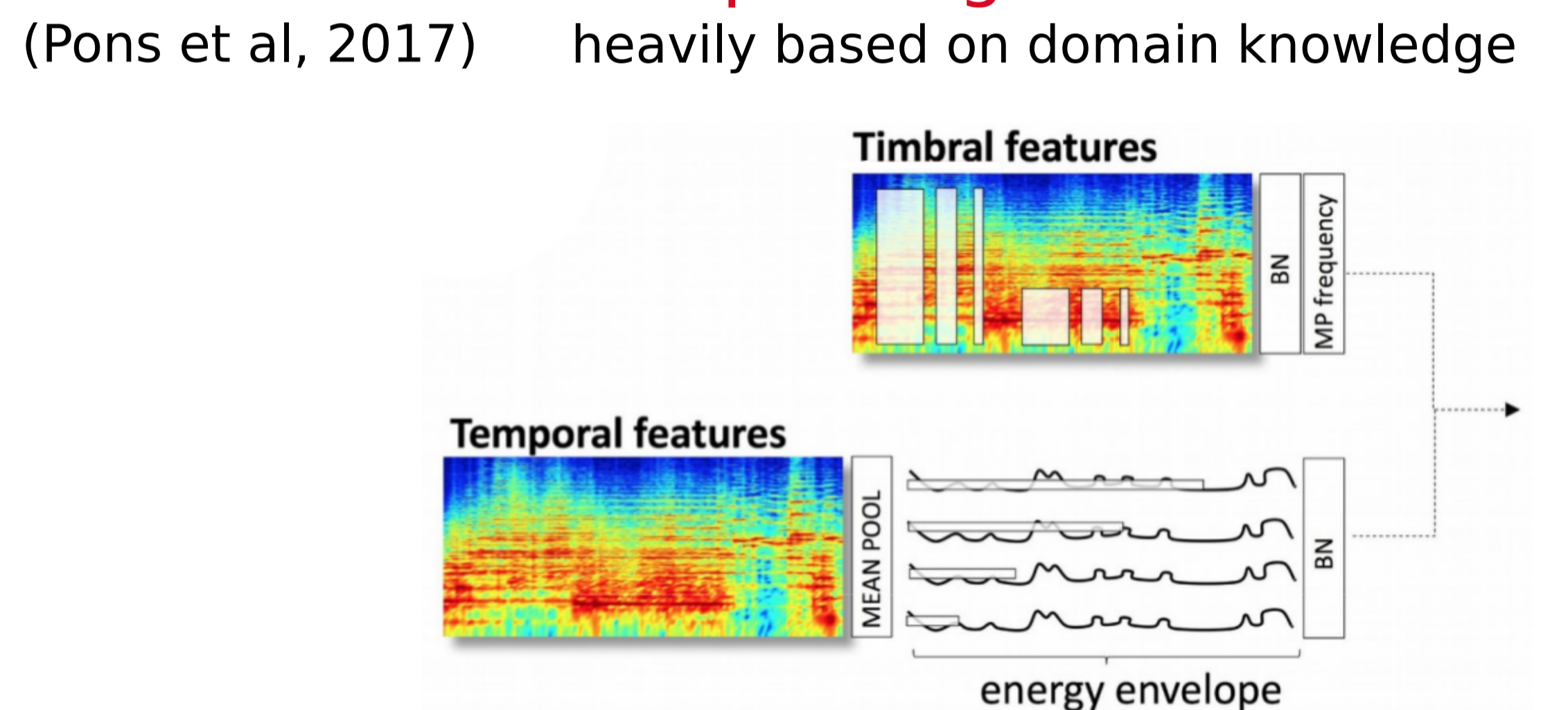
Many CNN frontends!

DESIGN BASED ON DOMAIN KNOWLEDGE?	FILTERS CONFIG?	INPUT SIGNAL? <i>waveform</i> end-to-end learning in the strictest sense	<i>pre-processed waveform</i> which is generally formatted in 2D i.e.: time-frequency representation
yes	single filter shape in 1st CNN layer	<b>FRAME-LEVEL</b>  filter length: 512 stride: 256 (Dieleman et al., 2014)	<b>VERTICAL OR HORIZONTAL</b>  filter shape: 7x90 (Lee et al., 2009) OR filter shape: 7x3 (Schlüter et al., 2014)
yes	many filter shapes in 1st CNN layer	<b>FRAME-LEVEL</b>  filter lengths: 512, 256, 128 stride: 64 (Zhu et al., 2016)	<b>VERTICAL AND/OR HORIZONTAL</b>  vertical filter shapes: 3x40, 1x75 horizontal filter shapes: 1x3, 1x10 (Pons et al., 2017)
no	minimal filter expression	<b>SAMPLE-LEVEL</b>  stack of 3x1 filters (Lee et al., 2017)	<b>SMALL RECTANGULAR FILTERS</b>  stack of 3x3 filters (Choi et al., 2016)

Waveform front-end  
an assumption-free model



Spectrogram front-end  
heavily based on domain knowledge



Waveform front-end:

- frame-level single-shape < frame-level many-shapes
- frame-level many-shapes << sample-level

Spectrogram front-end:

- domain knowledge intuitions are a valid guide for designing your front-end

Many possible backends!

Variable-length input back-end:

- max and average pooling
- attention models
- RNNs

**Note:** music is of variable length!

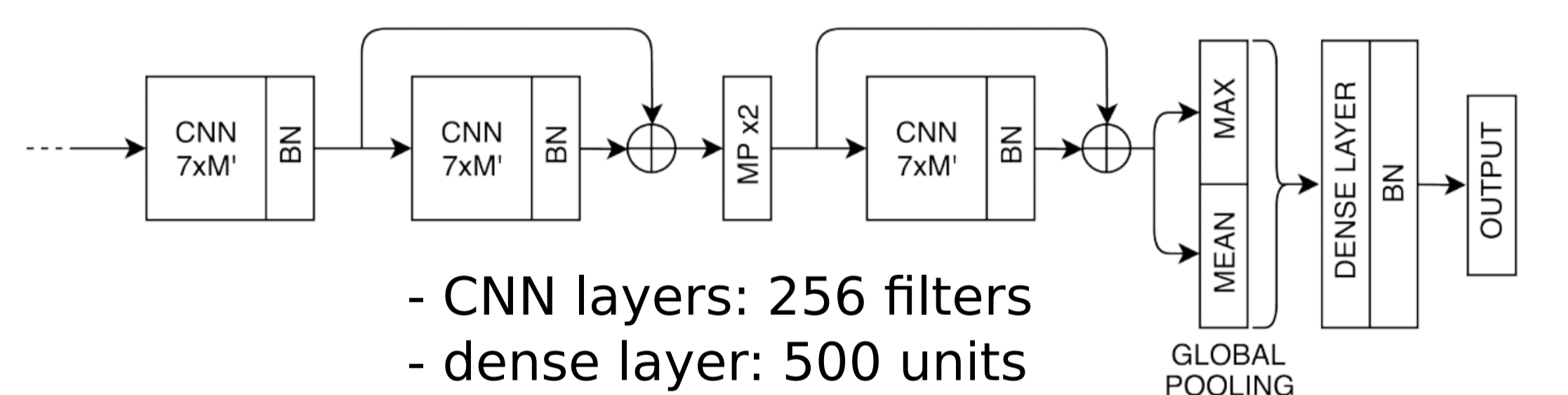
Fixed-length input back-end:

- fully convolutional models
- DNNs

**But:** most models assume a fixed-length input!

(Dieleman et al, 2014)

Shared back-end



## Quantitative results

Models	#training examples	ROC AUC	PR AUC	$\sqrt{MSE}$	$\Delta$ ROC AUC	$\Delta$ PR AUC	$\Delta$ $\sqrt{MSE}$	training time	Audio segments of 15 sec. Song-level predictions: - averaging windowed predictions
GBT+features	1.2M	91.61%	54.27%	0.1569	-	-	-	-	
Waveform	1M	91.54%	57.86%	0.1501	0.6%	1.49%	0.0021	< 2 weeks	
Spectrogram	1M	<b>92.14%</b>	<b>59.35%</b>	<b>0.1480</b>					
Waveform	500k	91.23%	56.15%	0.1537	0.54%	1.75%	0.0044	≈ 1 week	Annotations, 2 distributions: - bi-modal, classification tags ROC-AUC and PR-AUC
Spectrogram	500k	91.76%	57.90%	0.1493					- uniform, regression tags
Waveform	100k	89.16%	49.25%	0.1591	0.97%	2.83%	0.0049	few days	Error
Spectrogram	100k	90.13%	52.08%	0.1542					

## Qualitative results

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

Predicting each tag independently vs. predicting all tags together

"East Coast", "West Coast" / "baroque period", "classic period"

Reproduce this experiment online:

[jordipons.me/apps/music-audio-tagging-at-scale-demo](http://jordipons.me/apps/music-audio-tagging-at-scale-demo)

## Conclusions

Better performance than GBT+features baseline

spectrogram front-ends > waveform front-end  
..but the gap has been reduced!  
with more training data and Lee et al. front-end

Models' implementation in tensorflow:

[github.com/jordipons/music-audio-tagging-at-scale-models](https://github.com/jordipons/music-audio-tagging-at-scale-models)