

Language	Subjective preference (%) in naturalness					$p$ value
	LSTM	Concat	WaveNet (L)	WaveNet (L+F)	No preference	
North American English	23.3	<b>63.6</b>			13.1	$\ll 10^{-9}$
	18.7		<b>69.3</b>		12.0	$\ll 10^{-9}$
	7.6			<b>82.0</b>	10.4	$\ll 10^{-9}$
		32.4	<b>41.2</b>		26.4	0.003
		20.1		<b>49.3</b>	30.6	$\ll 10^{-9}$
			17.8	<b>37.9</b>	44.3	$\ll 10^{-9}$
Mandarin Chinese	<b>50.6</b>	15.6			33.8	$\ll 10^{-9}$
	25.0		23.3		51.8	0.476
	12.5			<b>29.3</b>	58.2	$\ll 10^{-9}$
		17.6	<b>43.1</b>		39.3	$\ll 10^{-9}$
		7.6		<b>55.9</b>	36.5	$\ll 10^{-9}$
			10.0	<b>25.5</b>	64.5	$\ll 10^{-9}$

Table 2: Subjective preference scores of speech samples between LSTM-RNN-based statistical parametric (**LSTM**), HMM-driven unit selection concatenative (**Concat**), and proposed WaveNet-based speech synthesizers. Each row of the table denotes scores of a paired comparison test between two synthesizers. Scores of the synthesizers which were significantly better than their competing ones at  $p < 0.01$  level were shown in the bold type. Note that **WaveNet** (L) and **WaveNet** (L+F) correspond to WaveNet conditioned on linguistic features only and that conditioned on both linguistic features and  $F_0$  values.