

# BWSNet: Automatic Perceptual Assessment of Audio Signals

Clement Le Moine Veillon\*<sup>1</sup>, Victor Rosi\*<sup>2</sup>, Pablo Arias Sarah<sup>3</sup>, Léane Salais<sup>1</sup>, Nicolas Obin<sup>1</sup>

<sup>1</sup> STMS Lab - IRCAM, CNRS, Sorbonne Université, Paris, France

<sup>2</sup> Department of Speech Hearing and Phonetic Science, University College London, London, UK

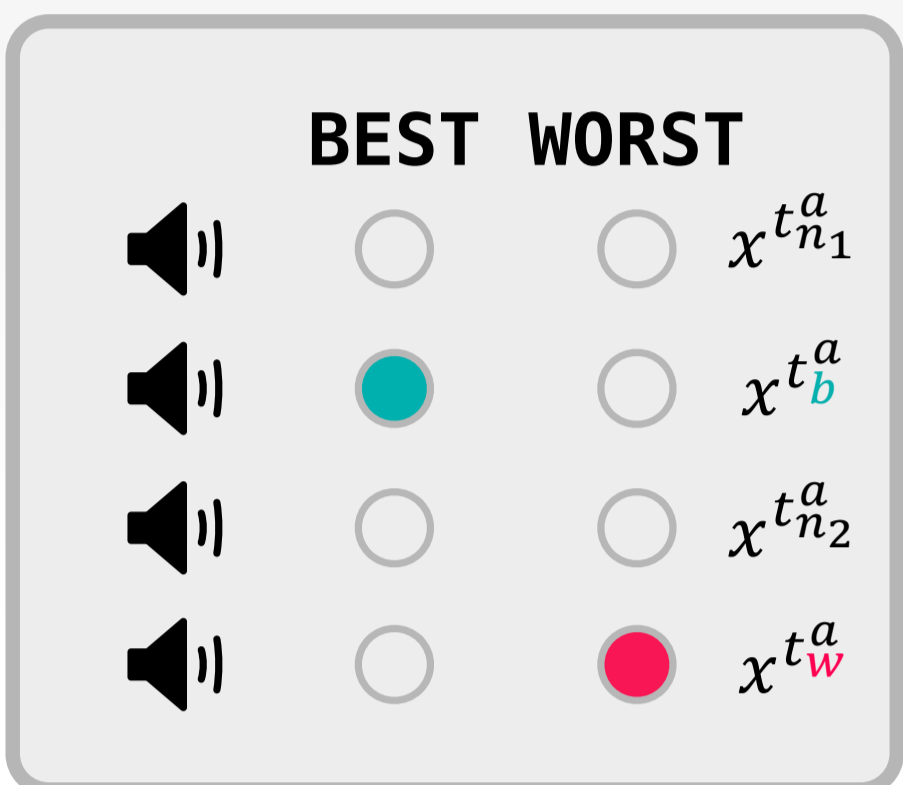
<sup>3</sup> School of Psychology and Neuroscience, University of Glasgow, Glasgow, UK

## Context

Revealing the **perceptual representation** of any sensory object usually requires conducting an experiment with human participants evaluating stimuli, based on a specific criterion. Recently, attention was given to **Best-Worst Scaling**<sup>1</sup> (BWS) as a method to assess perceptual qualities of sounds<sup>2,3</sup>. In this paper, we introduce **BWSNet**, a model for automatic perceptual assessment based on BWS data in a **metric learning** task.

## From BWS trial $t^a$ ...

« Which samples are perceived as **most** and **least** A? »



Derived ordinal relations<sup>4</sup>

$$\begin{cases} x^{t_b^a} >_a x^{t_w^a} \\ x^{t_b^a} >_a x^{t_{n_i}^a} \\ x^{t_{n_i}^a} >_a x^{t_w^a} \end{cases}$$

$x >_a y$ :  
x is more A than y

## ... To distances

$h^x$ : BWS embedding

Inequality (1)

$$\|h^{t_b^a} - h^{t_w^a}\| \geq \|h^{t_b^a} - h^{t_{n_i}^a}\|$$

Inequality (2)

$$\|h^{t_b^a} - h^{t_w^a}\| \geq \|h^{t_w^a} - h^{t_{n_i}^a}\|$$

## Metric learning losses and optimisations

• **Relative Contrastive loss with dynamic margins  $\mathcal{L}_{drc}^{t^a}$**

$$\mathcal{L}_{drc}^{t^a} = \frac{1}{n_v^{t^a}} \sum_{i=1}^{N-2} \max(\|h^{t_b^a} - h^{t_{n_i}^a}\| - \|h^{t_b^a} - h^{t_w^a}\| + \alpha_{b,n_i}, 0) + \frac{1}{n_v^{t^a}} \sum_{i=1}^{N-2} \max(\|h^{t_w^a} - h^{t_{n_i}^a}\| - \|h^{t_b^a} - h^{t_w^a}\| + \alpha_{w,n_i}, 0)$$

• **Dynamic margin constraint  $\mathcal{L}_{dmc}^{t^a}$**

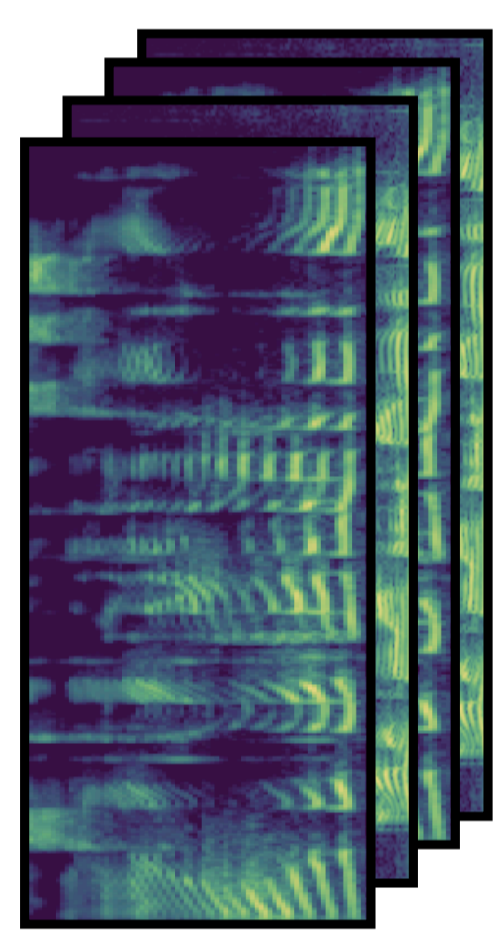
$$\mathcal{L}_{dmc}^{t^a} = \sum_{i=1}^{N-2} \gamma(\alpha_{b,n_i} - \mu) + \gamma(\alpha_{w,n_i} - \mu) \rightarrow \text{penalise the model's tendency to learn low-value margins.}$$

• **Fulfilled relations constraint  $\mathcal{L}_{frc}^{t^a}$**

$$\mathcal{L}_{frc}^{t^a} = \frac{n_v^{t^a}}{N} \rightarrow \text{Ensures that a decrease of } \mathcal{L}_{dmc}^{t^a} \text{ corresponds to more fulfilled relations in } t^a.$$

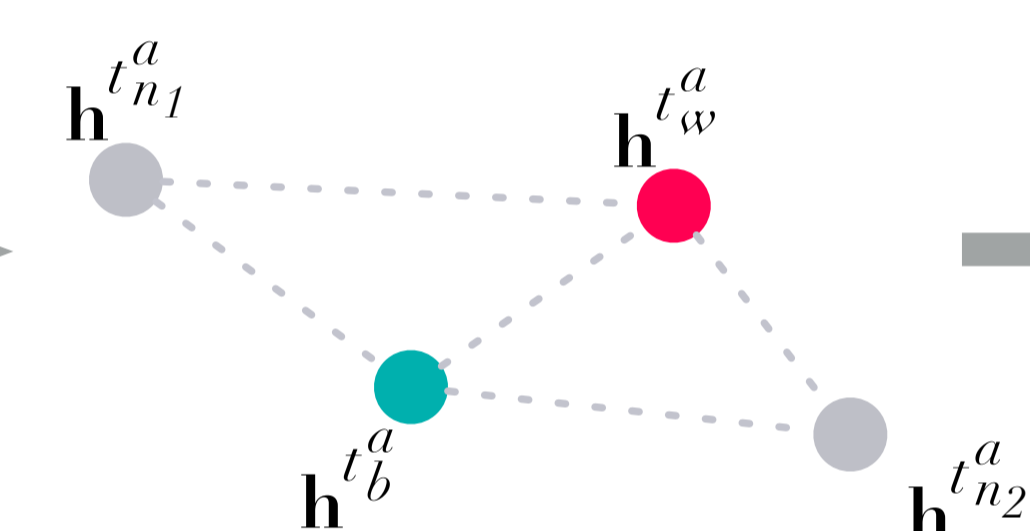
$$\text{Global BWSNet loss} \rightarrow \mathcal{L}^{t^a} = \mathcal{L}_{drc}^{t^a} + \lambda_{dmc} \mathcal{L}_{dmc}^{t^a} + \lambda_{frc} \mathcal{L}_{frc}^{t^a} \quad \lambda_{dmc}, \lambda_{frc} \geq 0$$

Mel Spectrograms of samples in trial  $t^a$



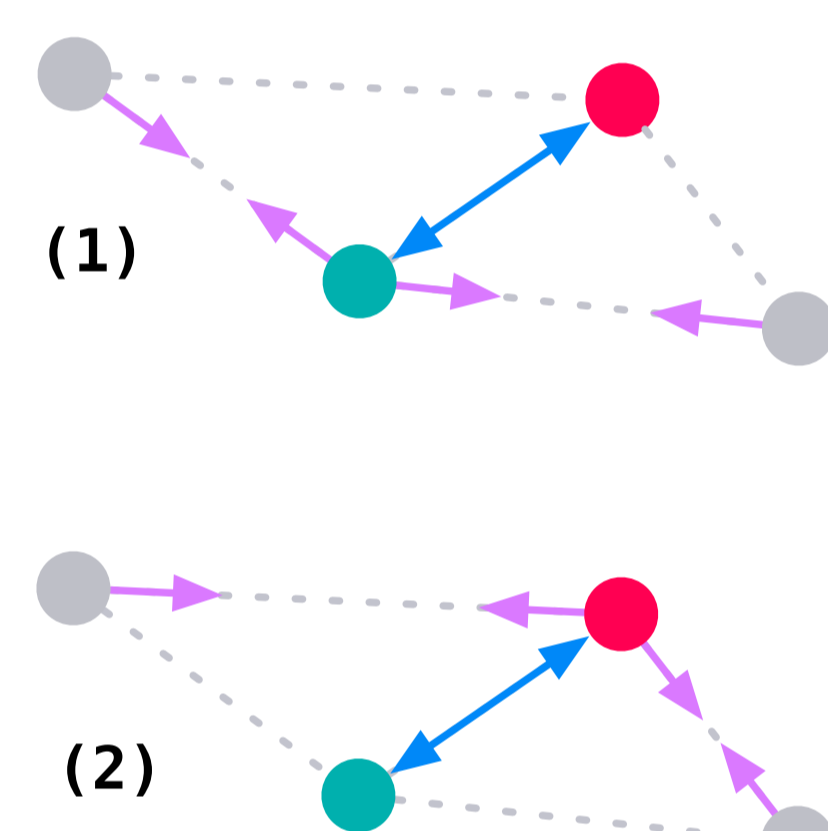
$\mathcal{L}^{t^a}$   
BWSNet training loss back-propagation.

Latent Space featuring BWS embeddings of samples in trial  $t^a$



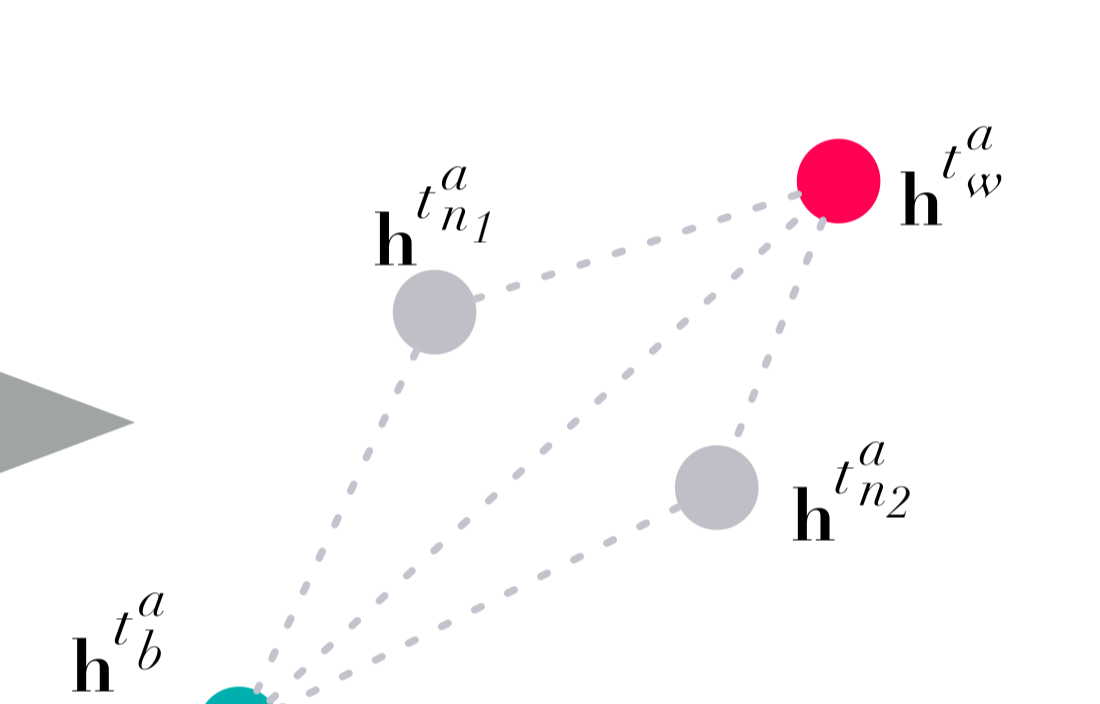
• Embedding for best  
• Embedding for worst  
• Embedding for neutrals

Training Phase



• reducing distance  
• expanding distance

Inference Phase



Latent space is structured according relations between sample from trial  $t^a$

## Application 1: Vocal Attitudes<sup>2</sup>

### Dataset

30-hour speech dataset showcasing four social attitudes: *friendliness*, *dominance*, *distance* and *seductiveness* in French.

### BWS experiment

96 participants evaluated four distinct subsets (N=2400) corresponding to each intended attitude.

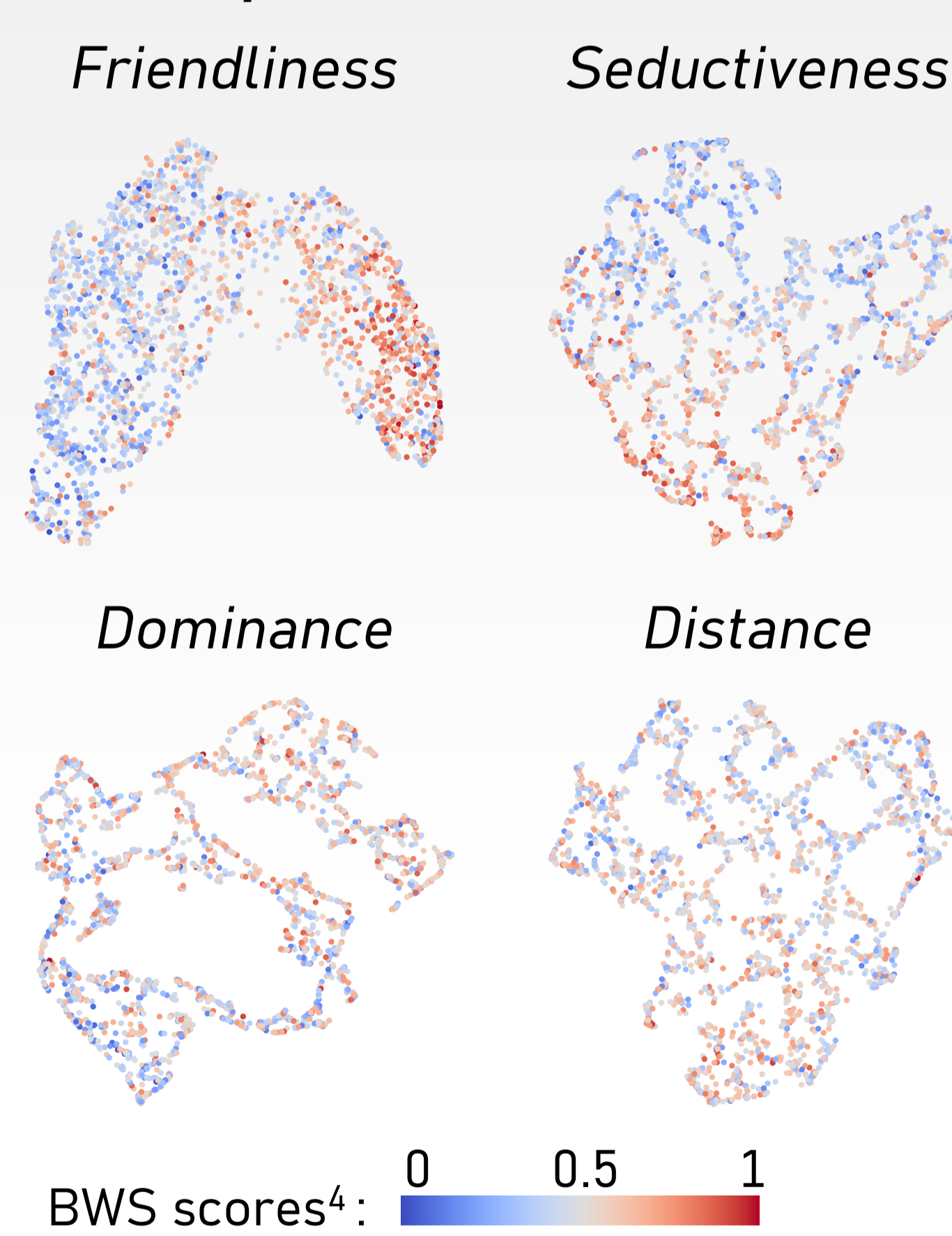
### Metric learning task

One perceptually structured latent space per attitude/dataset with BWSNet.

### Ablation study

Model	$\lambda_{dmc}$	$\lambda_{frc}$	FR* (%)
A-fixed	-	-	21.4 ± 8.5
A-l	0	0	1.0 ± 0.4
A-l-d	1	0	40.1 ± 18.7
A-l-d-fr	1	1	67.7 ± 4.5

### Latent space



BWS scores<sup>4</sup>: 0 0.5 1

## Application 2: Timbre Attributes<sup>3</sup>

### Dataset

Orchestral sounds dataset (N=520) showcasing multiple instruments and playing techniques.

### BWS experiment

16 sound experts evaluated the dataset on four timbral attributes: *brightness*, *warmth*, *roundness*, and *roughness*

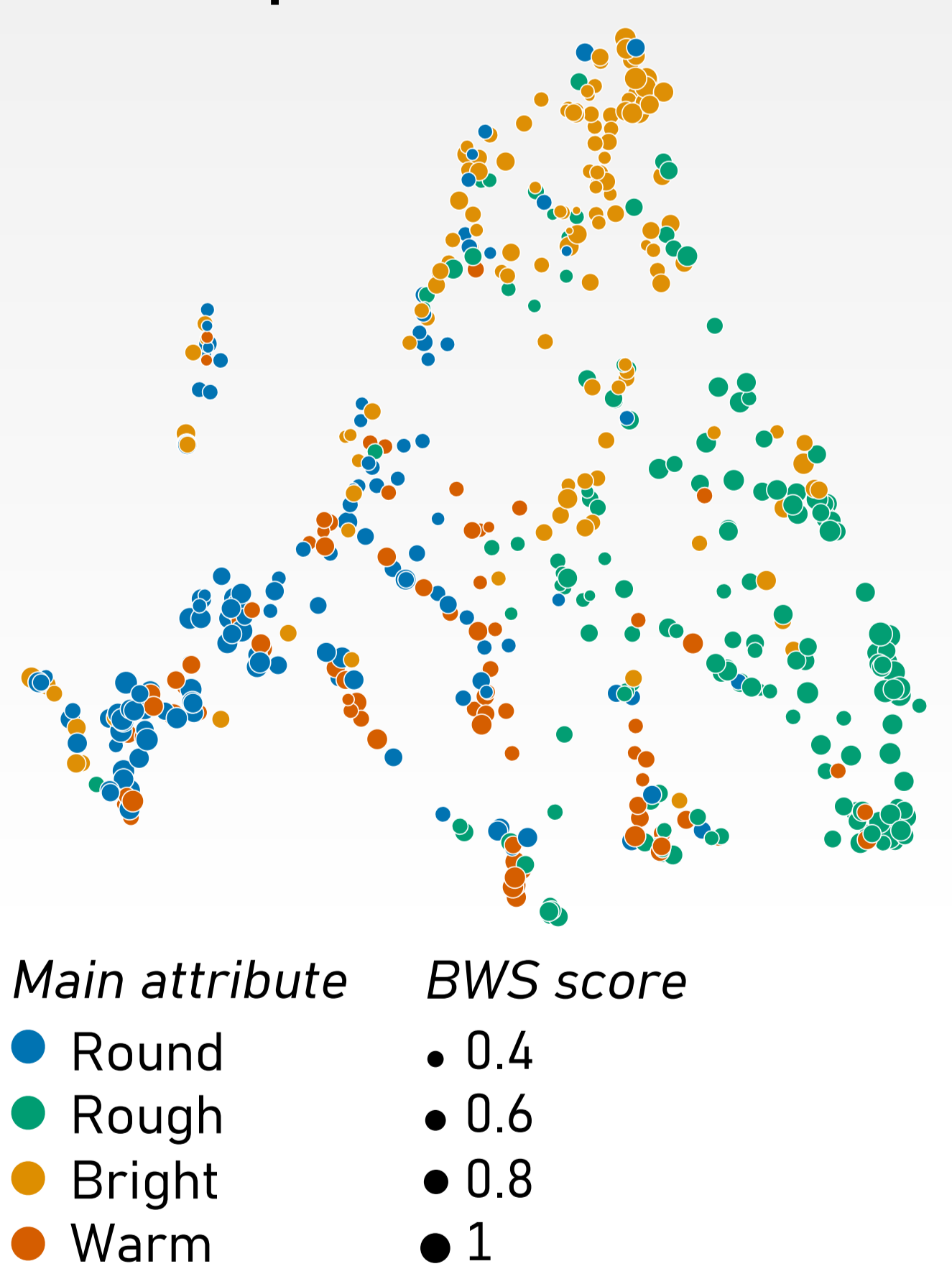
### Metric learning task

One perceptually structured latent space for all attributes with BWSNet.

### Ablation study

Model	$\lambda_{dmc}$	$\lambda_{frc}$	FR* (%)
A-fixed	-	-	37.5 ± 2.5
A-l	0	0	26.1 ± 2.1
A-l-d	1	0	51.6 ± 3.5
A-l-d-fr	1	1	56.3 ± 2.4

### Latent space



Main attribute BWS score  
• Round • 0.4  
• Rough • 0.6  
• Bright • 0.8  
• Warm • 1

## Conclusion

- BWSNet, a novel method designed for automated audio perceptual assessment based on human judgements.
- Perceptual organisation of latent spaces for vocal attitudes and timbral attributes.
- Perspectives of enhanced analyses of perceptual data by leveraging the dimensions of the latent space.
- Possible applications for conditioning synthesis/conversion models of speech and musical data.

## References

- <sup>1</sup> Louviere, J. J., Flynn, T. N., & Marley, A. A. J. (2015). *Best-worst scaling: Theory, methods and applications*. Cambridge University Press.
- <sup>2</sup> Le Moine, C. (2023). *Neural conversion of social attitudes in speech signals* (Doctoral dissertation, sorbonne université).
- <sup>3</sup> Rosi, V., Arias Sarah, P., Houix, O., Misdariis, N., & Susini, P. (2023). *Shared mental representations underlie metaphorical sound concepts*. Scientific Reports, 13(1), 5180.
- <sup>4</sup> Hollis, G., & Westbury, C. (2018). *When is best-worst best? A comparison of best-worst scaling, numeric estimation, and rating scales for collection of semantic norms*. Behavior research methods, 50, 115-133.