# Segment-level Metric Learning for Few-shot Bioacoustic Detection

**Detection and Classification of Acoustic Scenes and Events (DCASE) Workshop 2022**

Haohe Liu[1] , Xubo Liu[1] , Xinhao Mei[1] , Qiuqiang Kong[2] , Wenwu Wang[1] , Mark D. Plumbley[1]

[1]Centre for Vision, Speech and Signal Processing (CVSSP), University of Surrey, UK
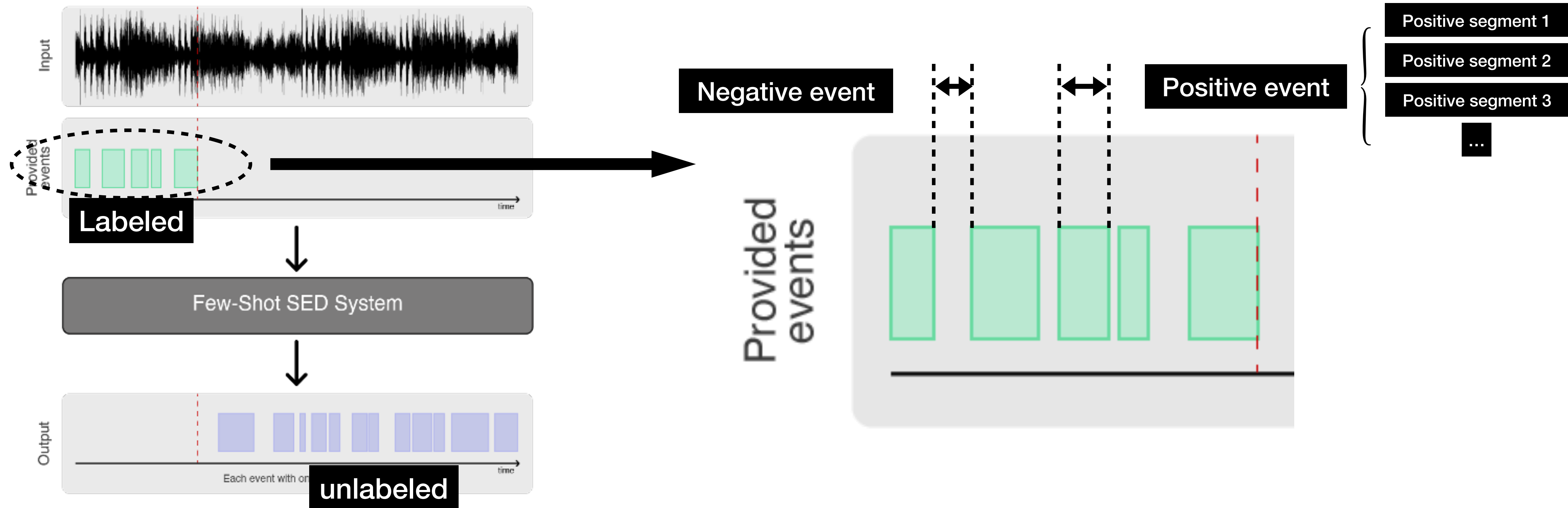[2]Speech, Audio, and Music Intelligence (SAMI) Group, ByteDance, China
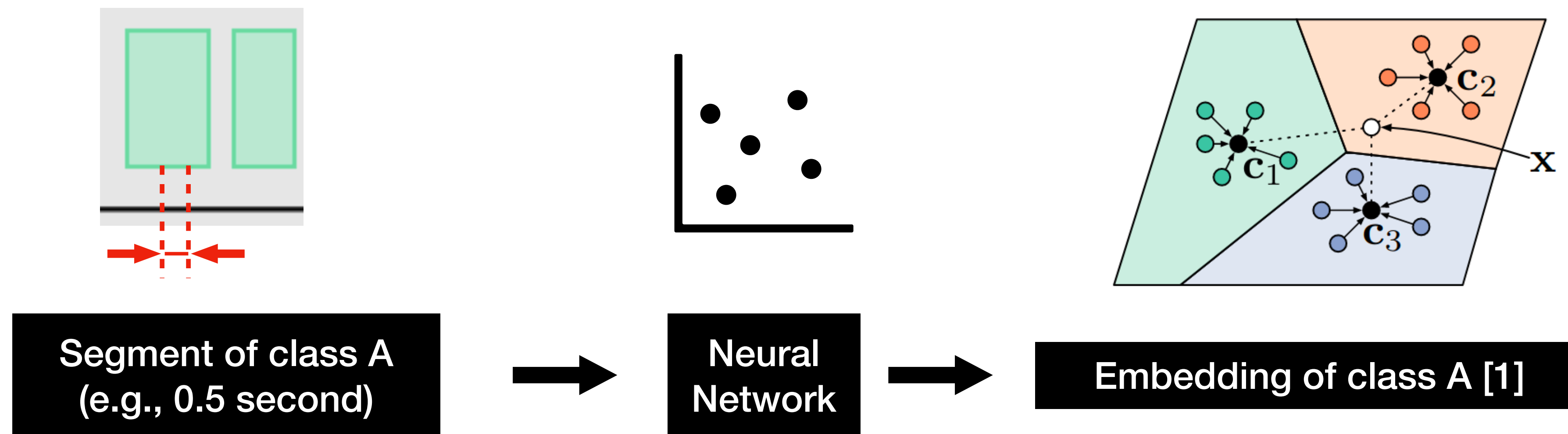
# Introduction
## Few-shot bioacoustic detection

- **Detecting** the occurrence time of **novel** sound events (e.g., new species) given a **few examples.**
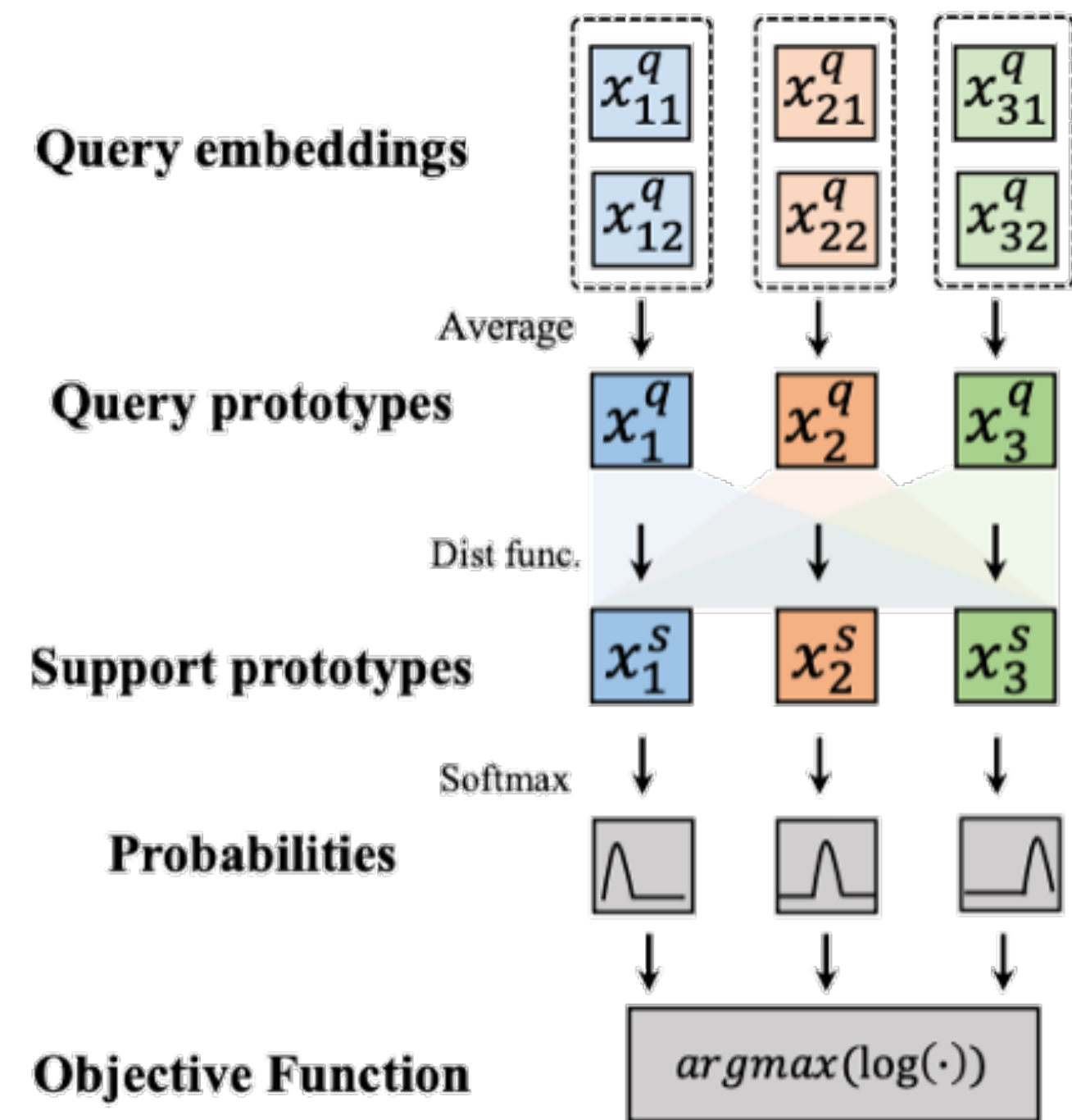


Picture taken from: https://github.com/c4dm/dcase-few-shot-bioacoustic

# Introduction
## Previous studies

- Metric learning → Prototypical network [1] → learn a latent space.

- In the latent space, the embeddings of different audio segments are expected to be

  - closer (the same class) or further apart (different classes).



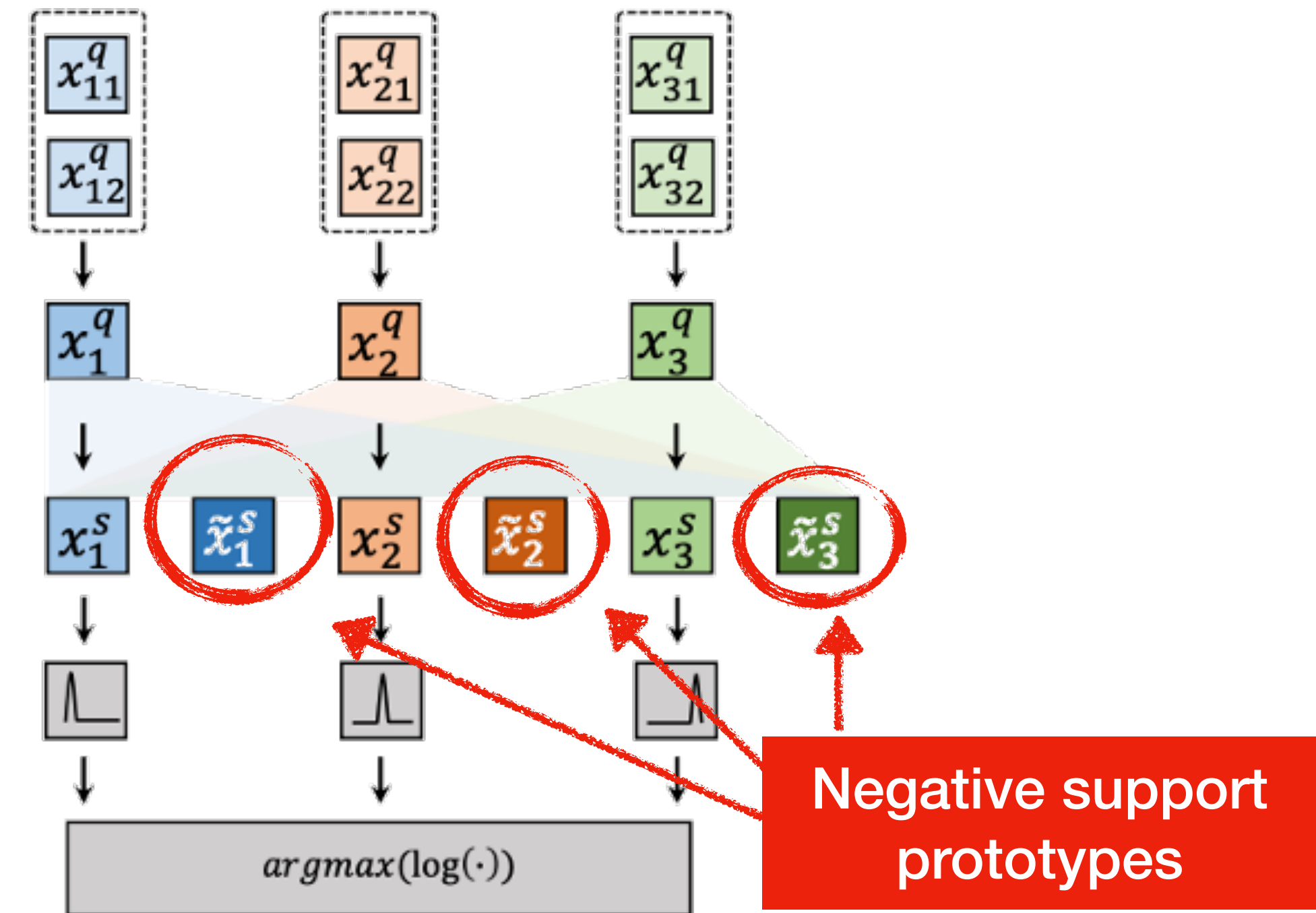**Segment of class A (e.g., 0.5 second)** ➡️ **Neural Network** ➡️ **Embedding of class A [1]**

[1] Snell, Jake, Kevin Swersky, and Richard Zemel. "Prototypical networks for few-shot learning." *Advances in neural information processing systems* 30 (2017).

# Proposed method



**Query embeddings**
$x_{11}^q$ $x_{21}^q$ $x_{31}^q$
$x_{12}^q$ $x_{22}^q$ $x_{32}^q$

Average

**Query prototypes** $x_1^q$ $x_2^q$ $x_3^q$

Dist func.

**Support prototypes** $x_1^s$ $x_2^s$ $x_3^s$

Softmax

**Probabilities**

**Objective Function** $argmax(\log(\cdot))$

Contrastive learning with negative segments

$x_{11}^q$ $x_{21}^q$ $x_{31}^q$
$x_{12}^q$ $x_{22}^q$ $x_{32}^q$

$x_1^q$ $x_2^q$ $x_3^q$

$x_1^s$ $\tilde{x}_1^s$ $x_2^s$ $\tilde{x}_2^s$ $x_3^s$ $\tilde{x}_3^s$

$argmax(\log(\cdot))$

Negative support prototypes

Previous studies

Training with
**positive** events
(<u>8.7%</u> of the training data)

Proposed learning with negative segments

Training with
**positive and negative** events
(<u>100%</u> of the training data)

# Proposed method
## Other Highlights

- **How to better adapt to the evaluation data?** → Transductive learning

- **Can more training data help?** → AudioSet strongly-labeled animal sound

- **Which feature is the best?** → Perform feature engineering.

- **Is F-score the ideal metric for evaluation?** → Evaluate with Polyphonic Sound Detection Score (PSDS).

- Other tricks applied: (i) Negative sample searching algorithm; (ii) Adaptive segment length; (iii) Augment training data; (iv) Post-processing.
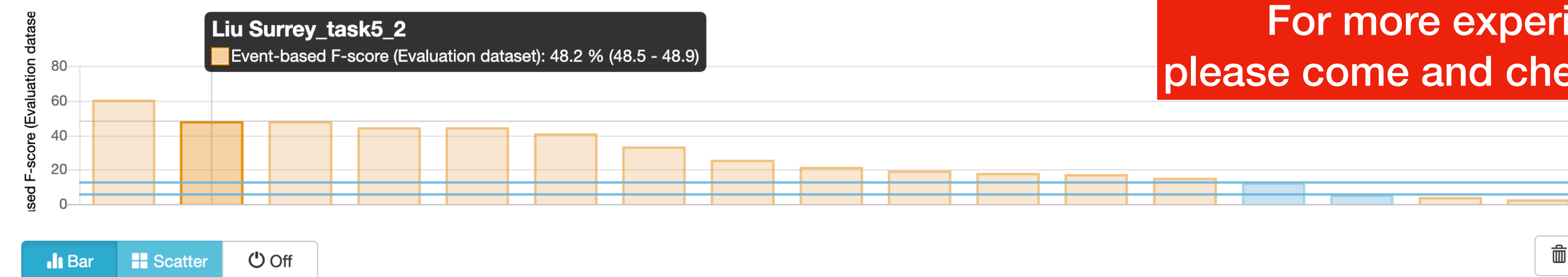
# Result

## DCASE2022-T5 Evaluation Set

- Based on the method proposed in this study, our system ranks 2$^{nd}$ in the DCASE 2022 Challenge Task 5: Few-shot Bio-acoustic Detection with an F-score of 48.2.



**For more experimental results please come and checkout our poster :-)**

# Thanks for your listening!

**In conclusion, the following points are helpful for few-shot bioacoustic detection:** (1) Contrastive learning with negative segments; (2) Transductive learning; (3) Data augmentations; (4) Feature engineering; and (5) External data from AudioSet.

- Paper: https://arxiv.org/abs/2207.07773

- Open-sourced code: https://github.com/haoheliu/DCASE_2022_Task_5