

Analyzing the effect of equal-angle spatial discretization on Sound event localization and detection



Saksham Singh Kushwaha^{1,2}, Iran R. Roman², Juan Pablo Bello^{2,3}

¹Courant Institute of Mathematical Sciences, New York University, NY, USA

²Music and Audio Research Lab, New York University, NY, USA

³Center for Urban Science and Progress, New York University, NY, USA



COURANT INSTITUTE
OF MATHEMATICAL SCIENCES



INTRODUCTION

Sound Event Localization and Detection (SELD)

- Detection of sound event occurrences with their spatial, temporal and class information.

Discretizing spatial targets for SELD

- Datasets typically use equal-angle discretization:
 $\theta = [-90, 90] \in \mathbb{Z}$ and
 $\phi = [-180, 180] \in \mathbb{Z}$ (in degrees).
- Problem:** this results in non-uniform spherical grid of points along the elevation axis
- An alternate discretization is Fibonacci [3], which results in a uniform spherical grid.

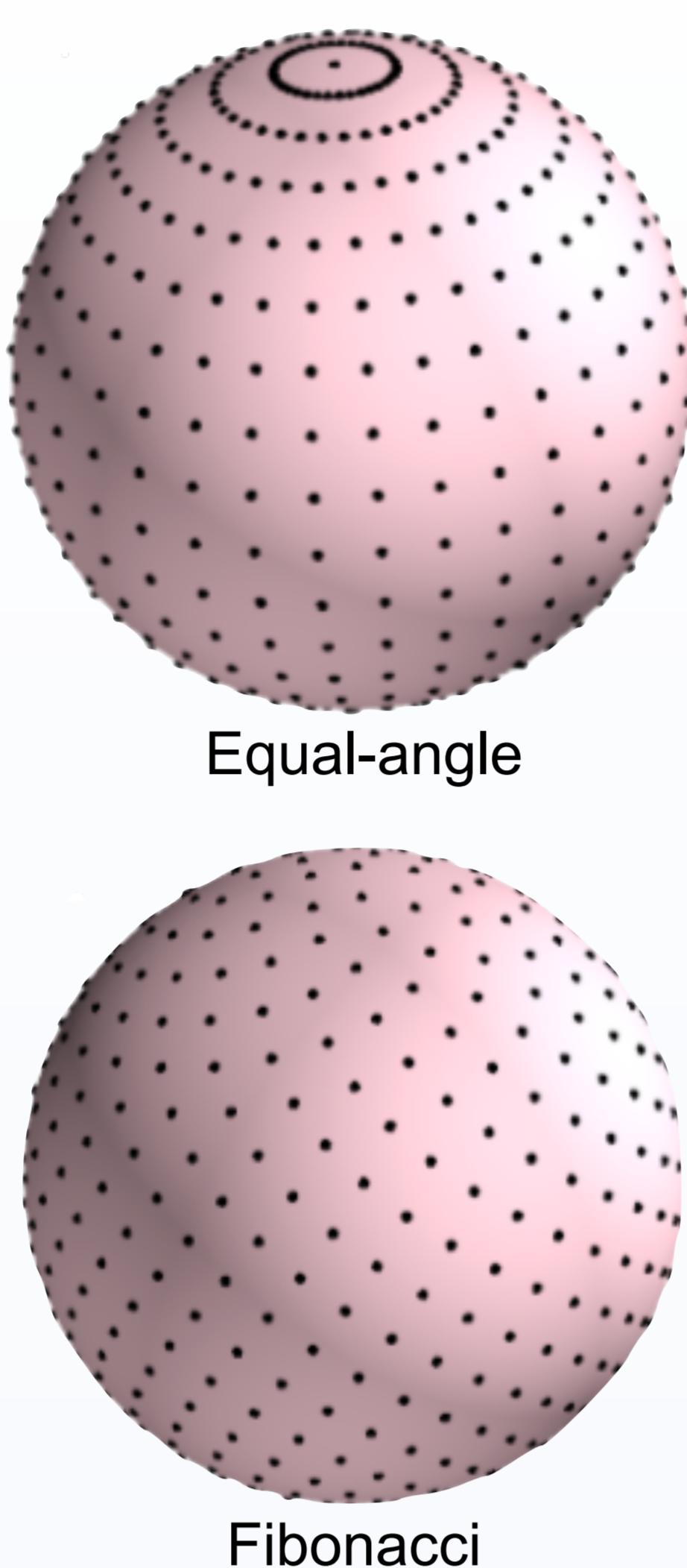
$$\mathbf{r}_n = [\cos(\phi_n) \sin(\theta_n), \sin(\phi_n) \sin(\theta_n), \cos(\theta_n)] \quad n = 1, \dots, N$$

where, $\phi_n = 2\pi n(1 - \frac{2}{1+\sqrt{5}})$ & $\theta_n = \arccos(1 - \frac{2n}{N})$

SELD Modeling

- DCASE-22 SELD baseline model is optimized with the ADPIT loss [2], which jointly minimizes detection and localization error using MSE.
- Hypothesis:** With equal-angle data, ADPIT will give larger localization errors at the equator.

Spatial discretization of sound event targets



Research questions

- Does equal-angle discretization affect model localization performance along the elevation axis?
- How can the model's localization performance be robust given the non-uniform grids used to discretize sound events in common SELD datasets?

ANALYSIS OF EQUAL-ANGLE VS UNIFORM DISCRETIZATION

Model

- SELD DCASE 2022 baseline model.

Data

- The DCASE STARSS2022[1] dataset features equal-angle targets (non-uniform grid along the elevation axis).

- In this dataset, more events occur around the equator.

- To control density of events, we generated uniformly-spatialized synthetic data along the elevation axis.

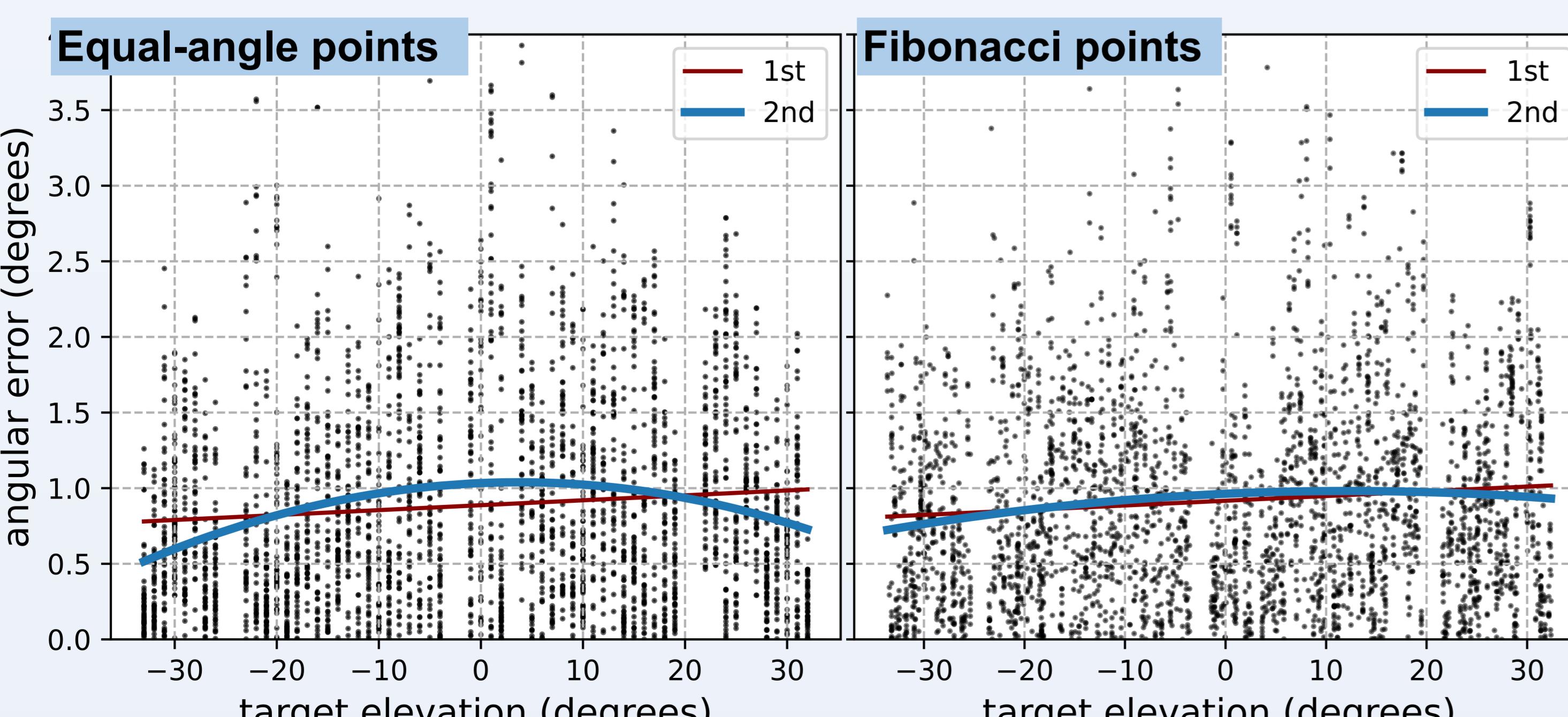
Experiments

- Model trained with targets that were discretized with either:
 - Equal-angle
 - Fibonacci

Results

- Equal-angle targets results in larger localization error around equator than pole
- Fibonacci targets show uniform localization error along the elevation axis

SELD angular localization error with:



Thresholded Angular error ADPIT (TAEADPIT) loss

The TAEADPIT Loss:

- Resampling with uniform discretization is not always feasible.
- TAEADPIT penalizes angular localization error, across azimuth and elevation.

$$l_{\alpha,ct}^{ALE} = \max(a_{\alpha,nct} ALE_{\alpha,nct}, H)$$

$$ALE_{\alpha,nct} = \angle(p(\mathbf{P}_{\alpha,nct}), p(\hat{\mathbf{P}}_{nct}))$$

$$L^{TAEADPIT} = \frac{1}{CT} \sum_c^C \sum_t^T \min_{\alpha \in \text{Perm}(ct)} l_{\alpha,ct}^{ACCDOATAE}$$

$$l_{\alpha,ct}^{ACCDOATAE} = \frac{1}{N} \sum_n^N MSE(\mathbf{P}_{\alpha,nct}^*, \hat{\mathbf{P}}_{nct}) + \beta(l_{\alpha,ct}^{ALE} - H)$$

Experiment 1

Research question

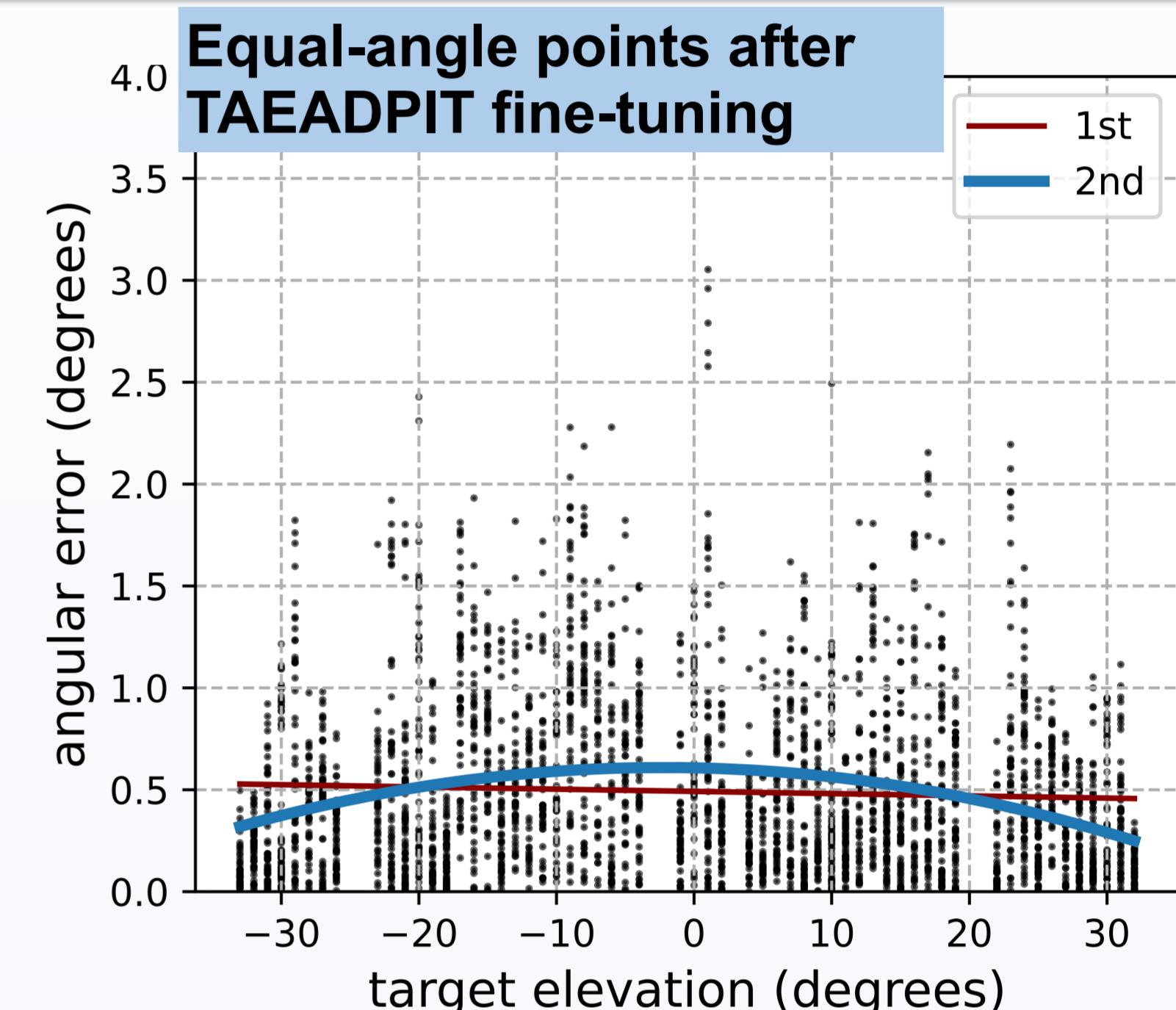
- Does fine-tuning with the TAEDPIT loss reduce equal-angle impact?

Methodology

- Use TAEADPIT to fine-tune the model previously trained with our synthetic equal-angle data.

Results

- Finetuning makes the localization error along the elevation axis more uniform.



Experiment 2

Research question

- Does using the TAEADPIT loss improve metrics on the DCASE 2022 SELD baseline model?

Methodology

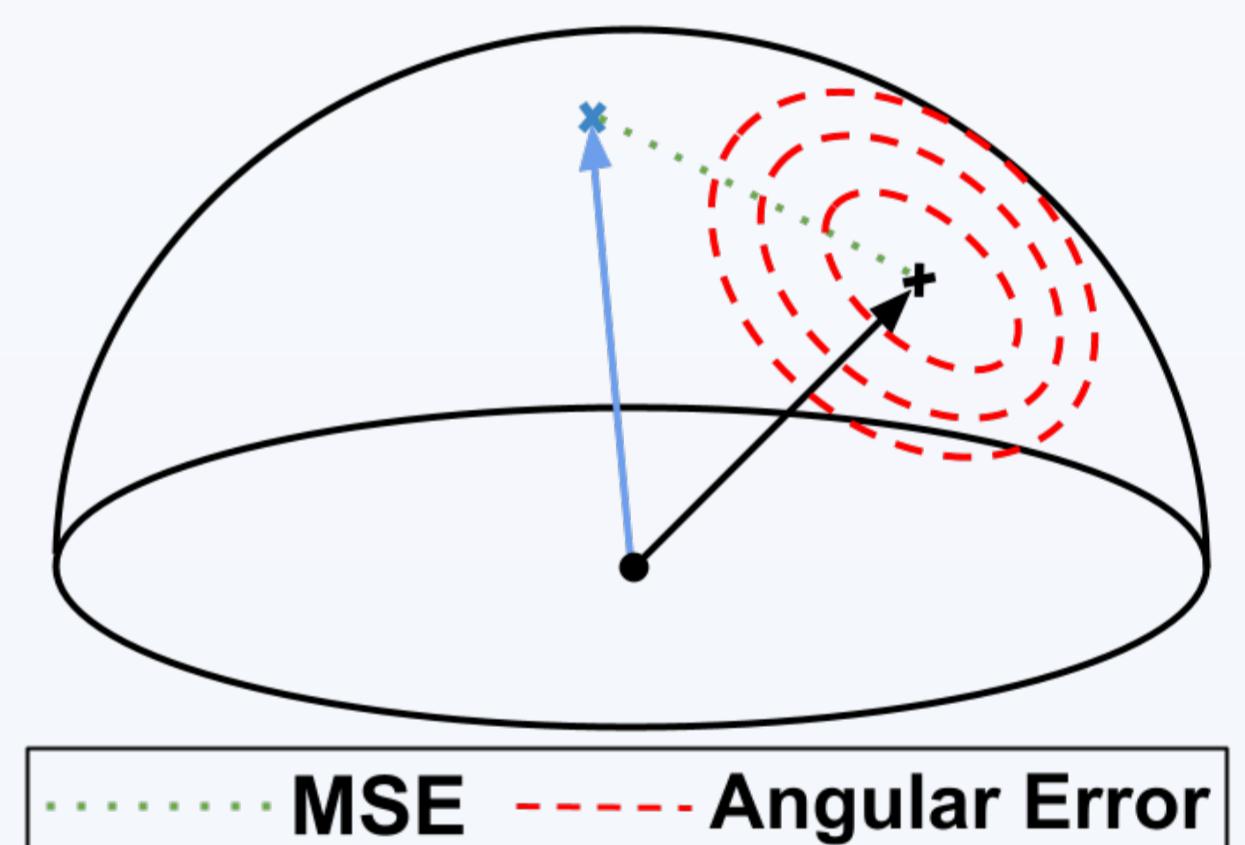
- DCASE22 SELD baseline model trained on STARS2022 (and supplementary synthetic) dataset with:
 - original ADPIT,
 - TAEDPIT fine-tune and from scratch
 - ADPIT with fibonacci targets.

Results

Loss	ER _{20°}	F _{20°}	LE _{CD}	LR _{CD}	SELD
ADPIT-base	0.69	0.24	30.43	0.43	0.55
TAEDPIT-tune	0.71	0.23	28.86	0.47	0.54
TAEDPIT	0.71	0.20	26.42	0.41	0.56
ADPIT-Fib	0.68	0.22	26.11	0.46	0.54

SELD model performance with different discretization and losses.
Training with TAEADPIT reduces localization error, similar to Fibonacci

Components of TAEADPIT loss



CONCLUSIONS AND FUTURE WORK

- The non-uniform localization error along the elevation axis can be mitigated by
 - Resampling with uniform spatial discretized targets
 - Finetuning with TAEADPIT
- In the future we would like to extend our study to other:
 - audio formats
 - target discretizations.

REFERENCES

- A. Politis, K. Shimada, P. Sudarsanam, S. Adavanne, D. Krause, Y. Koyama, N. Takahashi, S. Takahashi, Y. Mitsufuji, and T. Virtanen, "Starss22: A dataset of spatial recordings of real scenes with spatiotemporal annotations of sound events," *arXiv preprint arXiv:2206.01948*, 2022.
- K. Shimada, Y. Koyama, S. Takahashi, N. Takahashi, E. Tsunoo, and Y. Mitsufuji, "Multi-acccoa: Localizing and detecting overlapping sounds from the same class with auxiliary duplicating permutation invariant training," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Singapore, Singapore, May 2022.
- B. Keinert, M. Innmann, M. Saenger, and M. Stamminger, "Spherical fibonacci mapping," *ACM Transactions on Graphics (TOG)*, vol. 34, no. 6, pp. 1–7, 2015.