

Singer Identity Representation Learning Using Self-Supervised Techniques



Bernardo Torres¹, Stefan Lattner², Gael Richard¹

¹LTCI, Telecom Paris, Institut Polytechnique de Paris.

²Sony Computer Science Laboratories Paris



Introduction

Goal: obtain time-invariant identity representations from singing voice

Existing models from speech literature

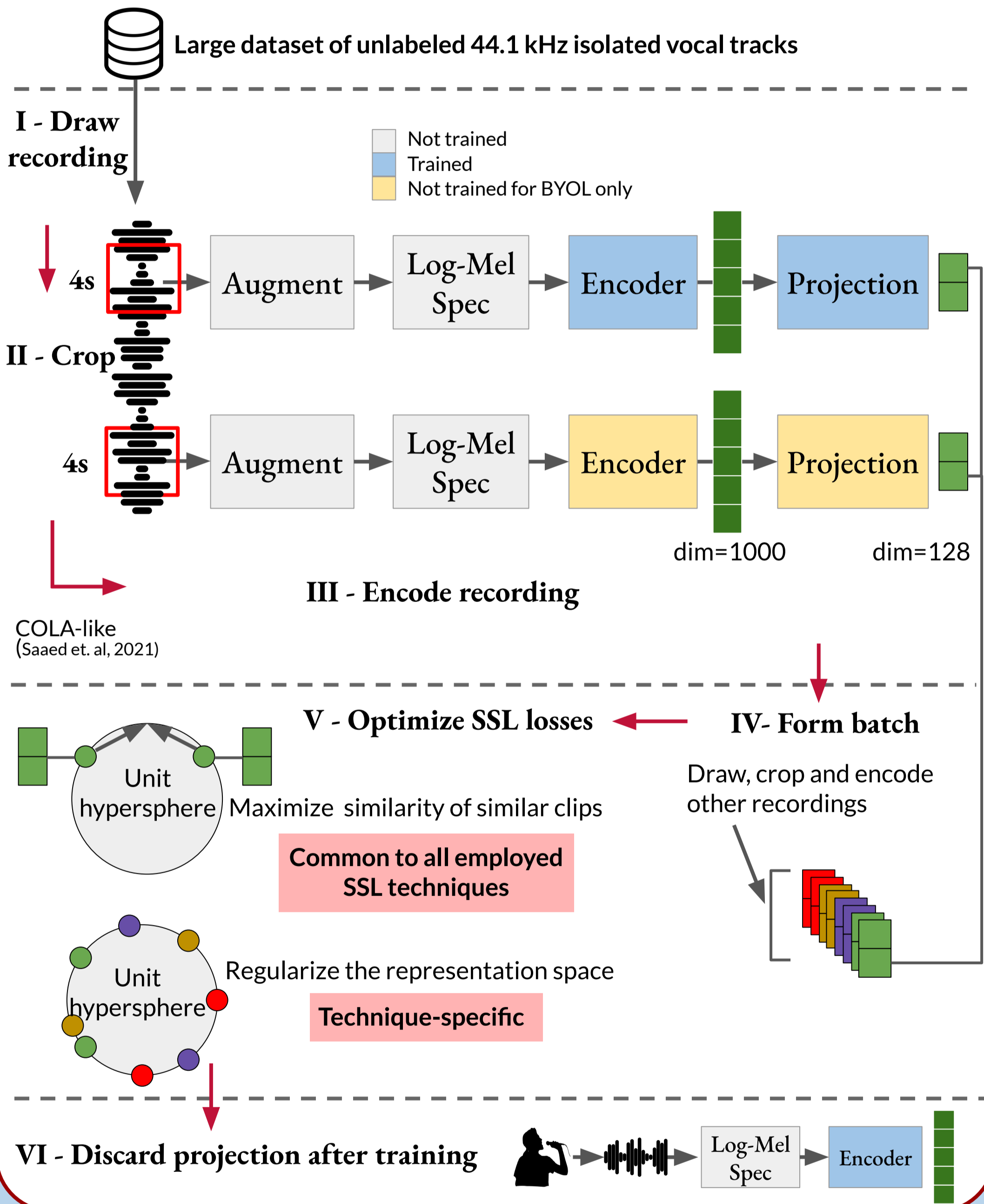
How well do models trained on speech generalize to singing voice?

Train identity extraction encoders

Lack of large labelled singing voice datasets

Can we train better models using Self-supervised Learning (SSL)?

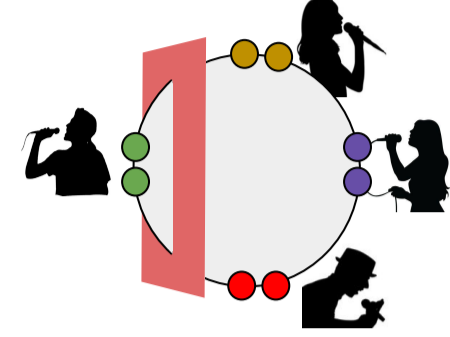
Overview and training



Evaluation

Singer identification

Linear classifier



Trained on embedding space (frozen encoder)
Test accuracy of N-fold cross validation

Singer similarity

Equal Error Rate (EER)

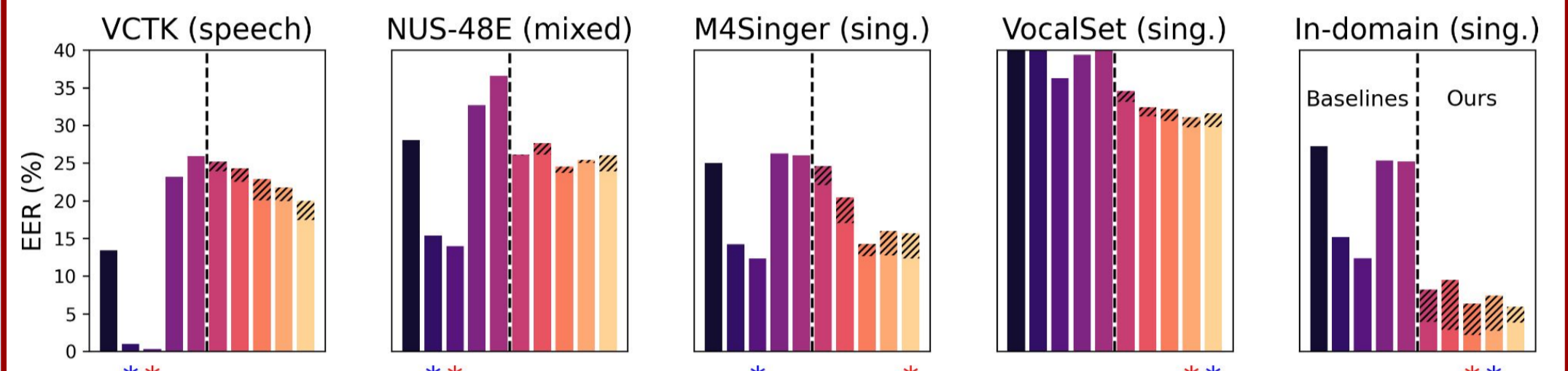
● open-set same/different binary classification

Mean Normalized Rank (MNR) Candidates

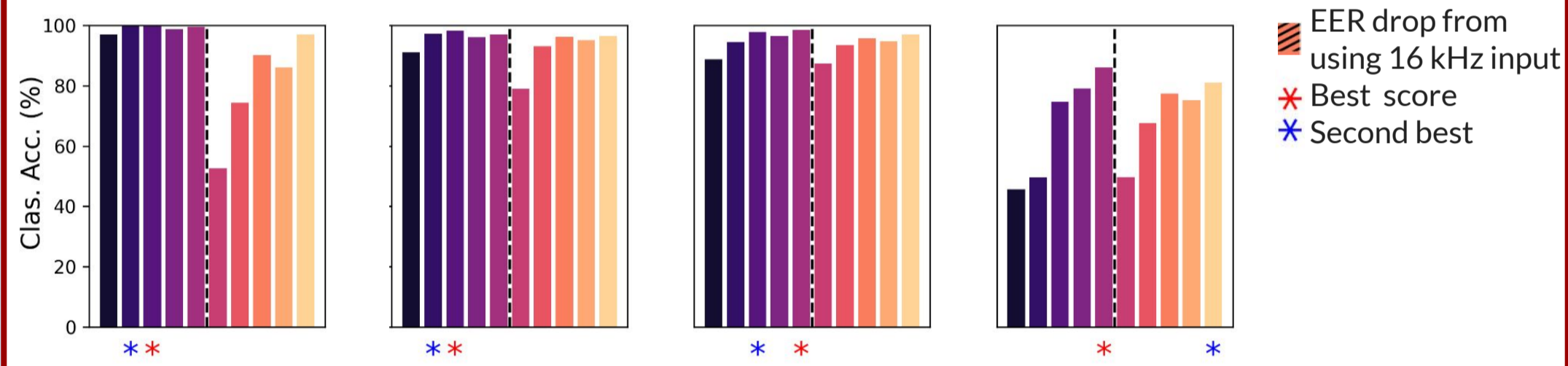
Rank ground-truth match Query by similarity with query

Results

Singer similarity (lower is better)



Singer identification (higher is better)



Summary of results for singing voice:

Baselines

Speech supervised

- GE2E
- F-ResNet
- H/ASP

↓ compared to evaluation on speech data

Work reasonably well; except for VocalSet

Speech SSL

- Wav2Vec-base
- XLSR-53

Bad on similarity, well on identification

Trained SSL models

- VICReg
- UNIF
- CONT
- CONT-VC
- BYOL

Comparable or superior to baselines

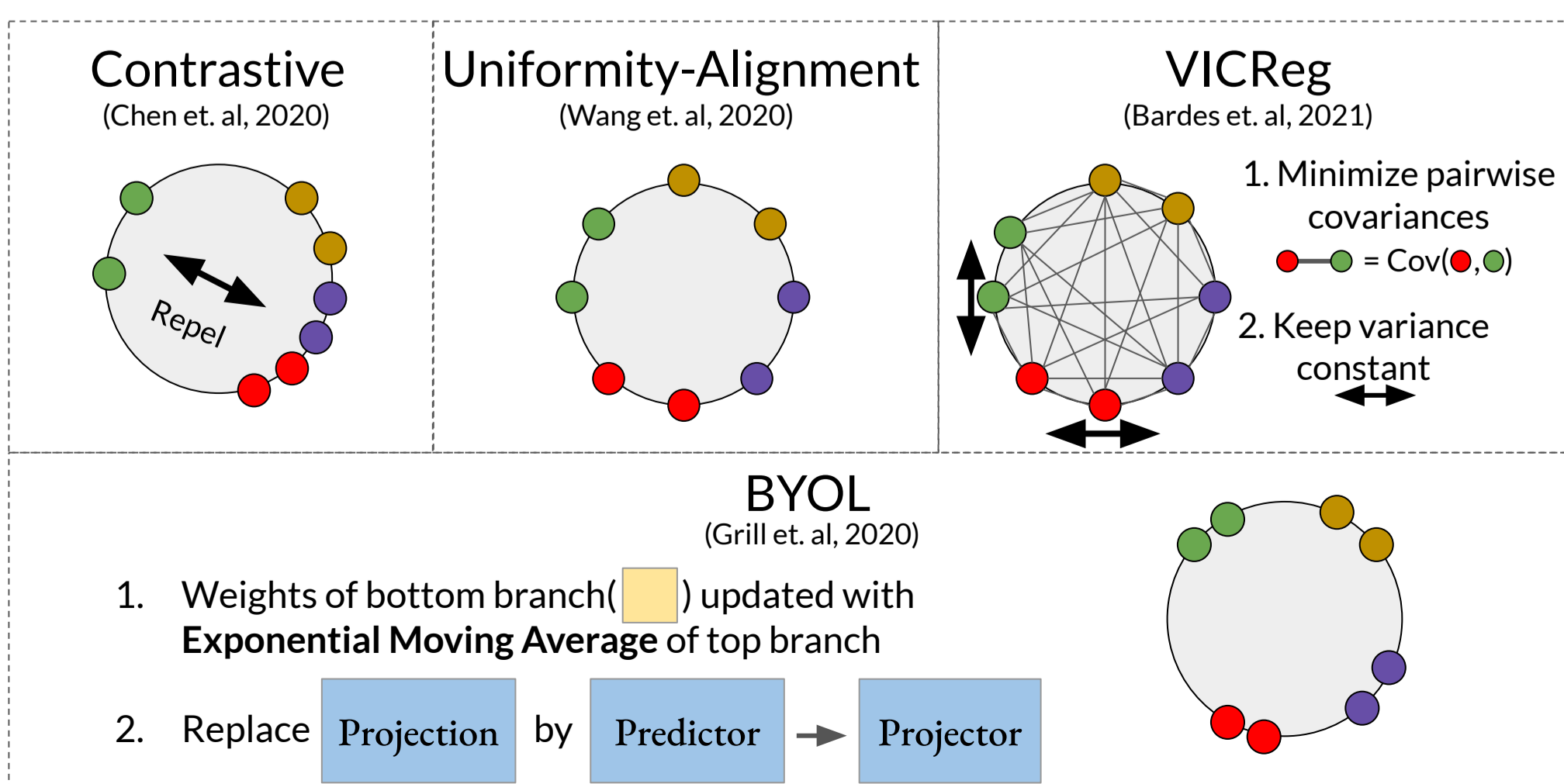
Best on out-of-domain: BYOL

Best In-domain: Contrastive

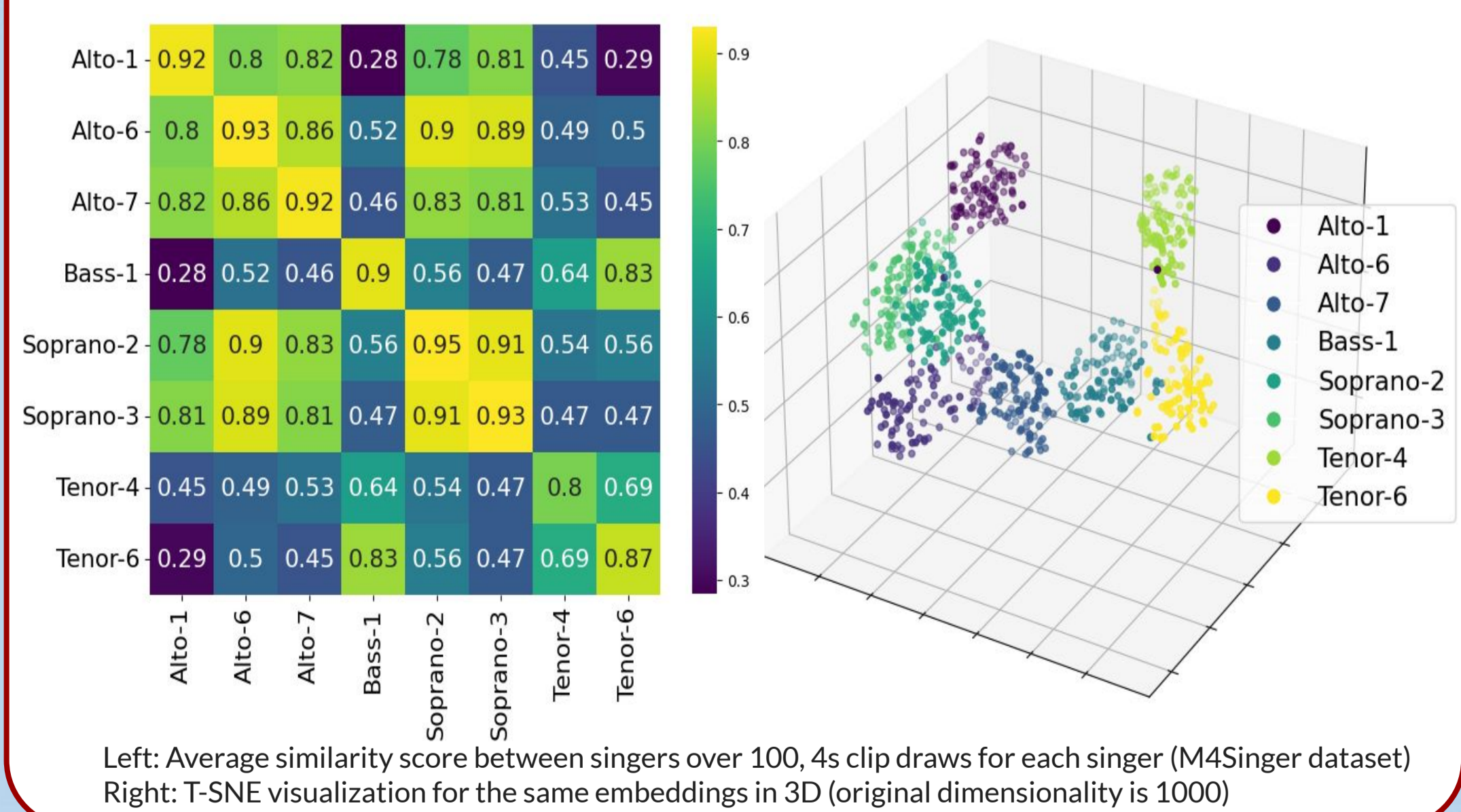
Self-supervised techniques

Common idea: representations from the same recording should be close

We trained models with the following SSL techniques:



Qualitative analysis



Conclusion

- Trained identity encoders using Self-Supervised Learning (SSL)
- Dataset: unlabeled isolated singing voice recordings
- Comparison with publicly available pre-trained speech models
- Evaluation on singer identification and similarity tasks
- A big gap still exists for challenging datasets
- Release of code and trained models

