

Neural Audio Fingerprint for High-specific Audio Retrieval based on Contrastive Learning

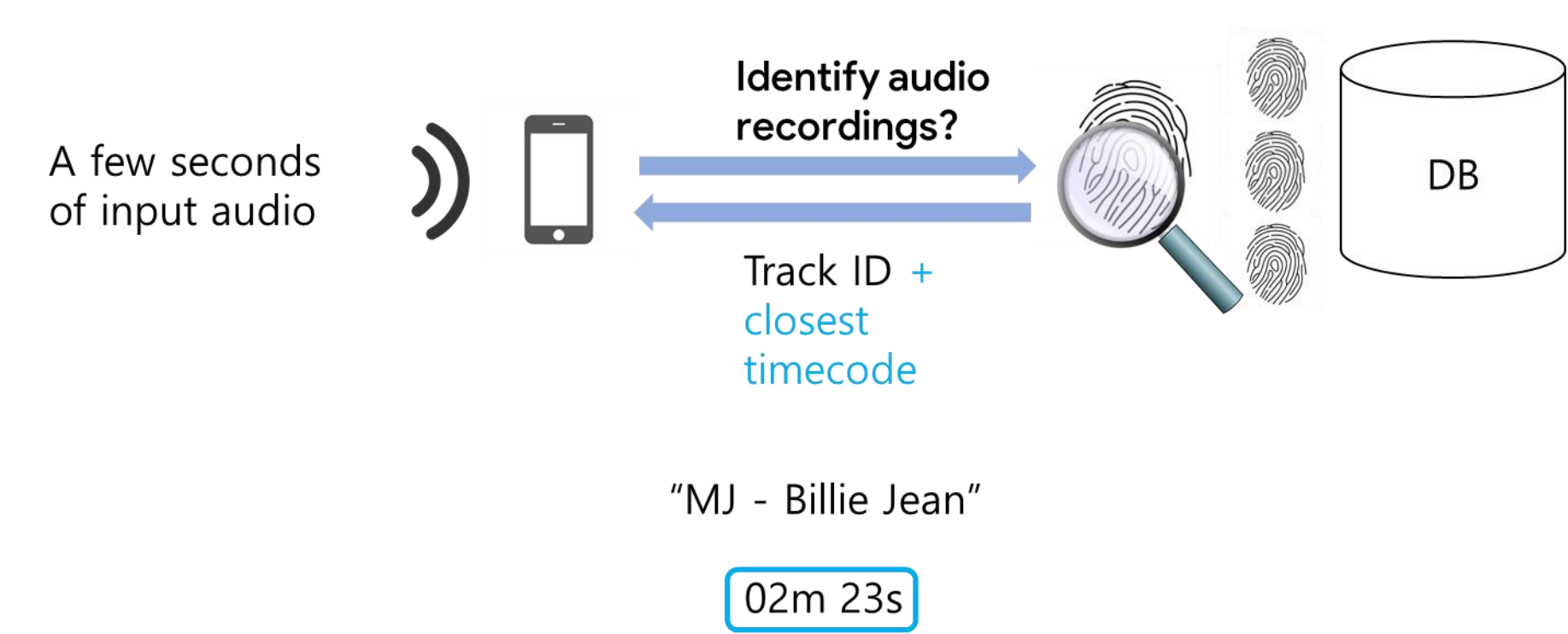
Sungkyun Chang¹ Donmoon Lee^{1,2} Jeongsoo Park¹ Hyungui Lim¹ Kyogu Lee² Karam Ko³ Yoonchang Han¹
Cochlear.ai, ²Seoul National University, ³SK Telecom

Abstract

Most of existing audio fingerprinting systems have limitations to be used for high-specific audio retrieval at scale. In this work, we generate a low-dimensional representation from a short unit segment of audio, and couple this fingerprint with a fast maximum inner-product search. To this end, we present a contrastive learning framework that derives from the segment-level search objective. Each update in training uses a batch consisting of a set of pseudo labels, randomly selected original samples, and their augmented replicas. These replicas can simulate the degrading effects on original audio signals by applying small time offsets and various types of distortions, such as background noise and room/microphone impulse responses. In the segment-level search task, where the conventional audio fingerprinting systems used to fail, our system using 10x smaller storage has shown promising results. Our code and dataset are available at <https://mimbres.github.io/neural-audio-fp/>.

Introduction

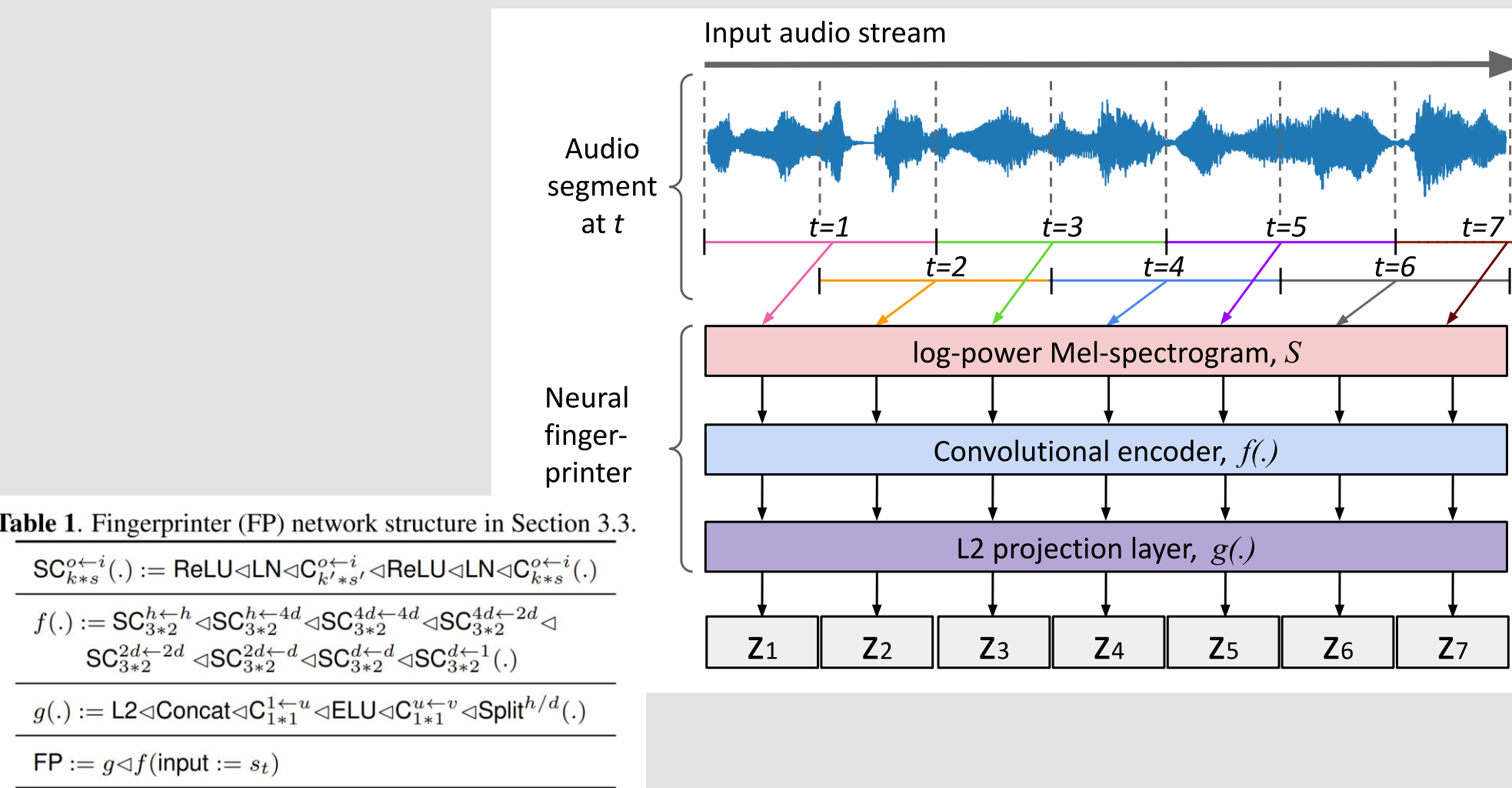
Goal: high-specific audio retrieval



Key aspects

- Segment-level audio search
 - allowing **miss-match less than 250 ms**
- Contrastive learning framework
 - simulating **maximum inner-product search (MIPS)** in mini-batch
 - explicitly sampling positive pairs to have **original-replica** relations
- Audio data augmentation chain
 - generating various types of **acoustic distractors** that helps training a robust audio fingerprint

Neural Audio Fingerprinter



- We employ $g \circ f : S \mapsto \mathcal{Z}(\cdot)$ as a segment-wise fingerprinter. It can generate fingerprint z_t that can represent a unit segment of 1 s audio x_t at the time step t .
- $f(\cdot)$ is a base encoder with separable convolution (SC) that computes internal representation.
- $g(\cdot)$ is a multi-head linear projection layer with L2 normalization.

Training $g \circ f(\cdot)$ with contrastive loss can be viewed as a common form of self-supervised learning (SSL). We maintain the self-supervised $g(\cdot)$ up to the final target task.

Contrastive Learning Framework

Due to the L2 normalization layer of $g(\cdot)$, we can use the inner-product $z_a^T z_b$ as a measure of similarity (z_a, z_b). Searching the most similar point (*) of database $V = \{v_i\}$ for a given query q in \mathcal{Z}^d space can be formulated as maximum inner-product search (MIPS):

$$v_i^{(*)} := \operatorname{argmax}_i(q^T v_i)$$

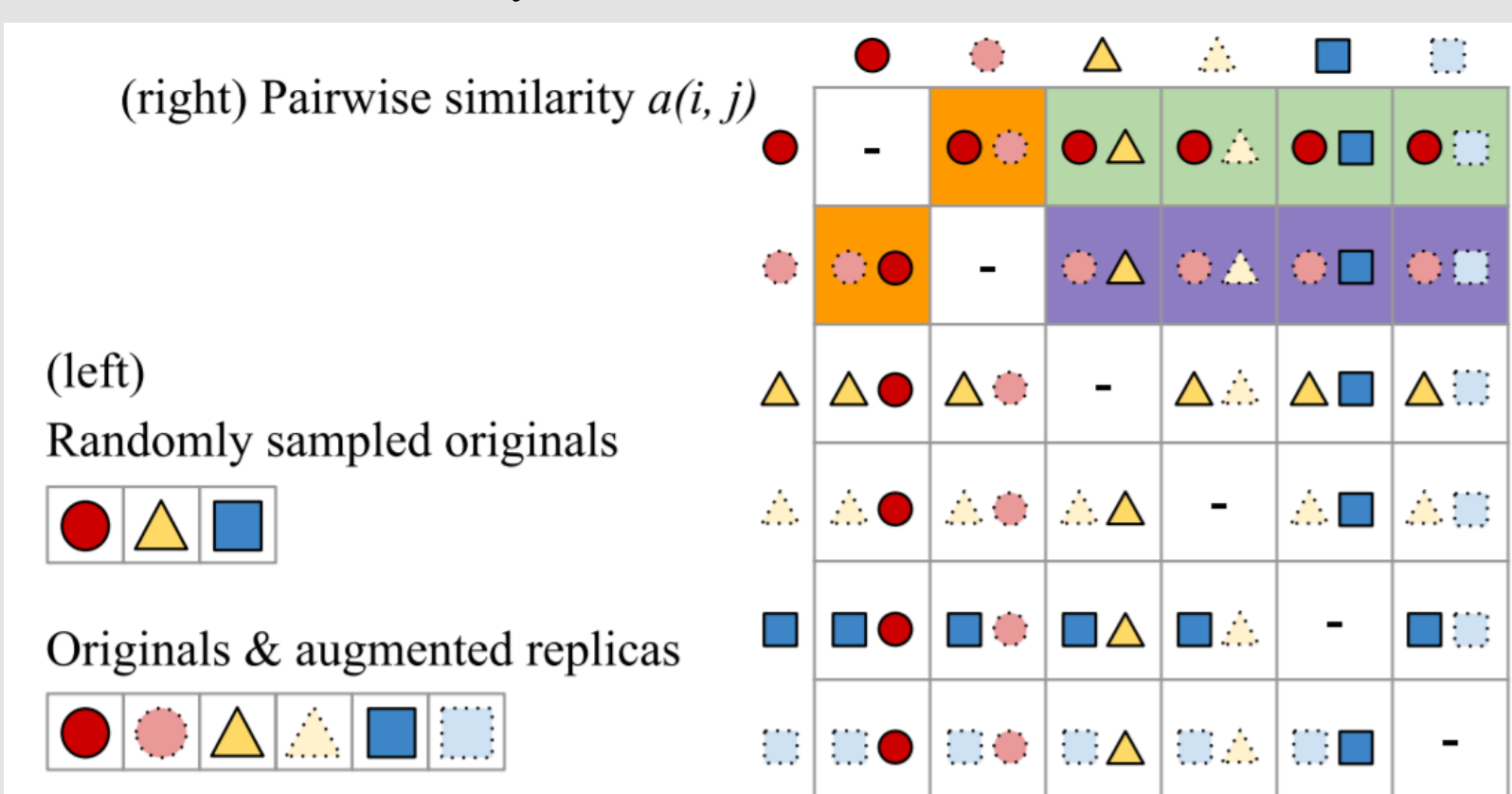


Fig. 2. Illustration of the contrastive prediction task in Section 2.1. (left) Batch size $N = 6$. We prepare $N/2$ pairs of original/replica. The same shapes with solid/dashed lines represent the positive pair of original/replica, respectively. (right) Each element in the matrix represents pairwise similarity. In each row, a prediction task can be defined as classifying a positive pair (one of the orange squares) against the negative pairs (green or purple squares) in the same row.

$$\ell(i, j) = -\log \frac{\exp(a_{i,j}/\tau)}{\sum_{k=1}^N \mathbf{1}(k \neq i) \exp(a_{i,k}/\tau)} \quad (1)$$

$\mathbf{1}(\cdot) \in \{0, 1\}$ is an indicator function that returns 1 iff (\cdot) is true, and $\tau > 0$ denotes the temperature parameter for softmax.

The Eq(1) can replace MIPS from the property: computing the top-k predictions in the softmax function is equivalent to the MIPS.

Experiments

- Dataset (from FMA)
 - Train 10K songs (30s cut each)
 - Test-DB: unseen 100K songs (full)
- Test-Query (synthesized)
 - 2K seq. generated with unseen augmentation source
- Augmentation source dataset (from AudioSet & others)
 - BG mix: 6.6 h environmental noise (subway, metro, pub, café,...) with “no music”
 - MIC and room/space IR
- Top-1 hit rate (%) = $100 \times \frac{(n \text{ of hits @Top-1})}{(n \text{ of hits @Top-1}) + (n \text{ of miss @Top-1})}$

Table 3. Top-1 hit rate (%) of large-scale (total of 100K songs) segment-level search. d denotes the dimension of fingerprint embedding. *exact match* means that our system finds the exact index. *near match* means a mismatch within ± 1 index or ± 500 ms.

Method	d	match	Query length in seconds					
			1 s	2 s	3 s	5 s	6 s	10 s
Now-playing (replicated)	128	exact	-	44.3	60.1	73.6	81.0	86.1
		near	-	46.8	63.5	75.2	81.6	86.3
Now-playing (modified for 1 s unit)	64	exact	25.8	58.5	69.3	78.5	81.4	87.7
		near	30.9	61.3	71.2	79.5	82.2	88.3
This work (N=640)	128	exact	26.3	58.2	69.5	78.4	81.4	87.8
		near	30.9	61.1	71.8	79.8	83.0	89.2
This work (N=320)	64	exact	54.6	78.9	85.4	90.4	92.0	94.9
		near	61.3	81.7	86.7	90.9	92.7	95.1
This work (N=120)	128	exact	61.0	82.2	87.1	91.8	93.1	95.2
		near	67.1	84.1	88.1	92.5	93.9	95.5
This work (no aug.)	128	exact	55.9	78.8	84.9	90.9	92.2	95.3
		near	62.3	80.9	86.3	91.5	92.8	95.5
This work (no aug.)	128	exact	0.0	0.0	0.0	0.0	0.0	0.0
		near	0.0	0.0	0.0	0.0	0.0	0.0

Table 4. Effect of fingerprint dimension d in 1 s segment search.

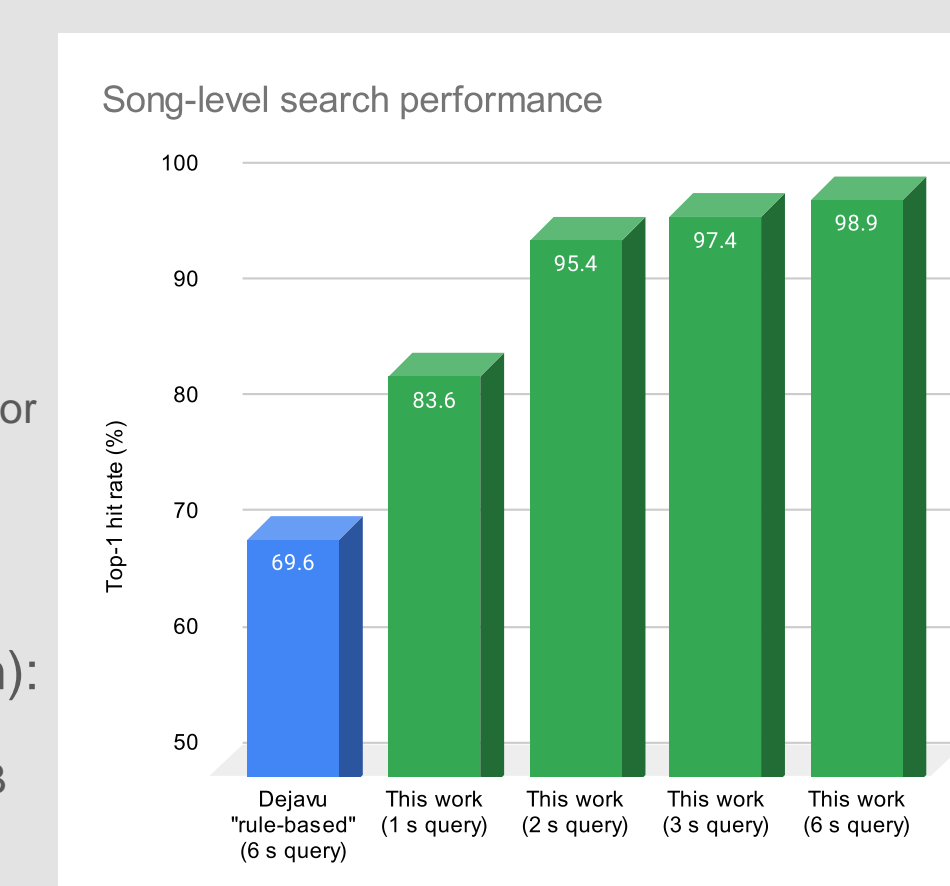
Embedding dimension	$d=16$	$d=32$	$d=64$	$d=128$
Top-1 hit rate @ 1 s (%)	11.6	40.2	54.6	62.2

- Now-playing < Now-playing mod.
- In every cases, our model outperformed over Now-playing. Contrastive learning framework > semi-hard triplet embedding.
- Augmentation was critical in training.
- The larger the batch size, the better the performance in all experiments.
- The longer the query sequence, the better the performance.
- Increasing embedding dimension d improves performance.

VS. Dejavu (rule-based)

In the small (10K-30s) DB test:
 - This work (6 s query) = 98.9 % exact match for segment-level search
 - Dejavu (6 s query) = 69.6% exact match for song-level search

Memory usage (w/ in-memory search):
 - Dejavu = 400 MB
 - This work ($d=64$, after compression) < 40 MB
 - This work uses 10x less memory!

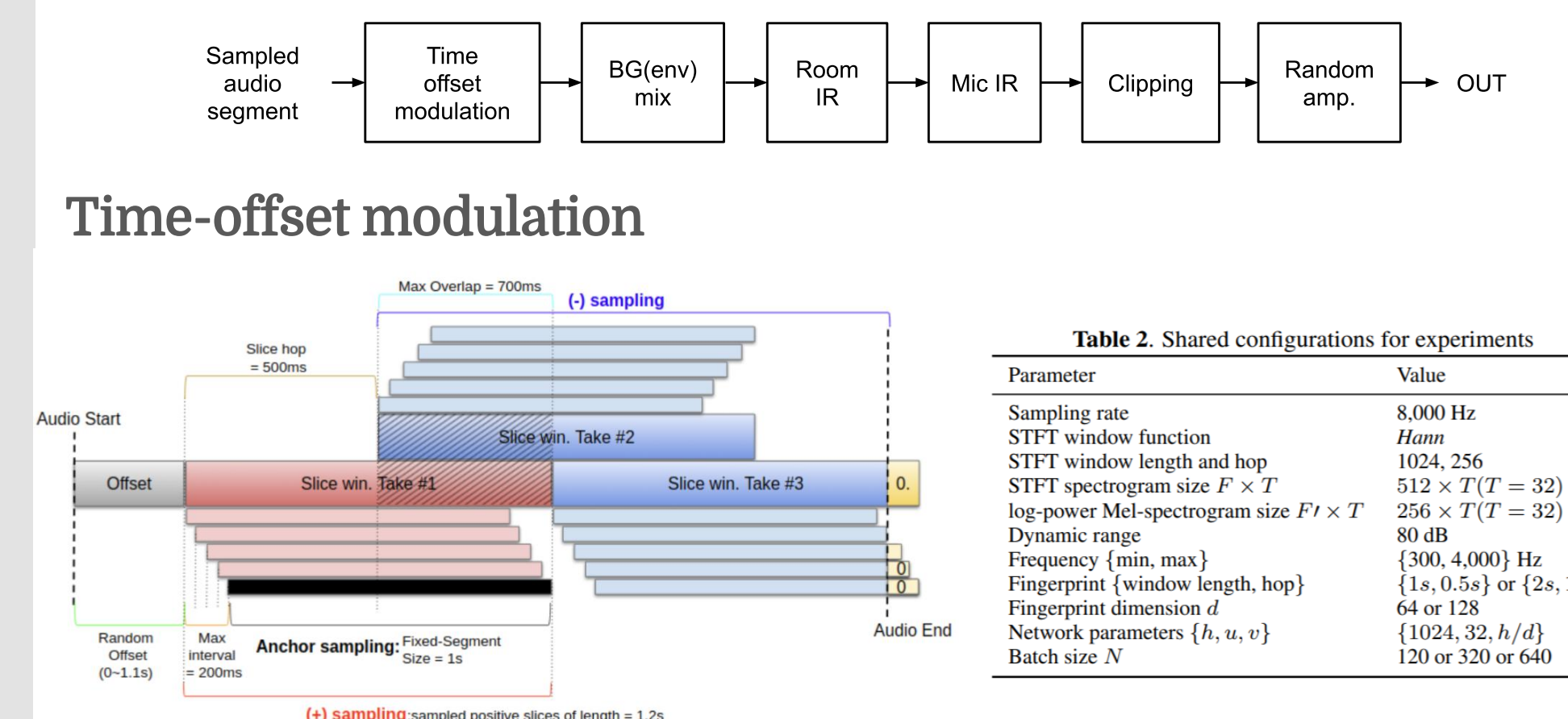


Conclusion

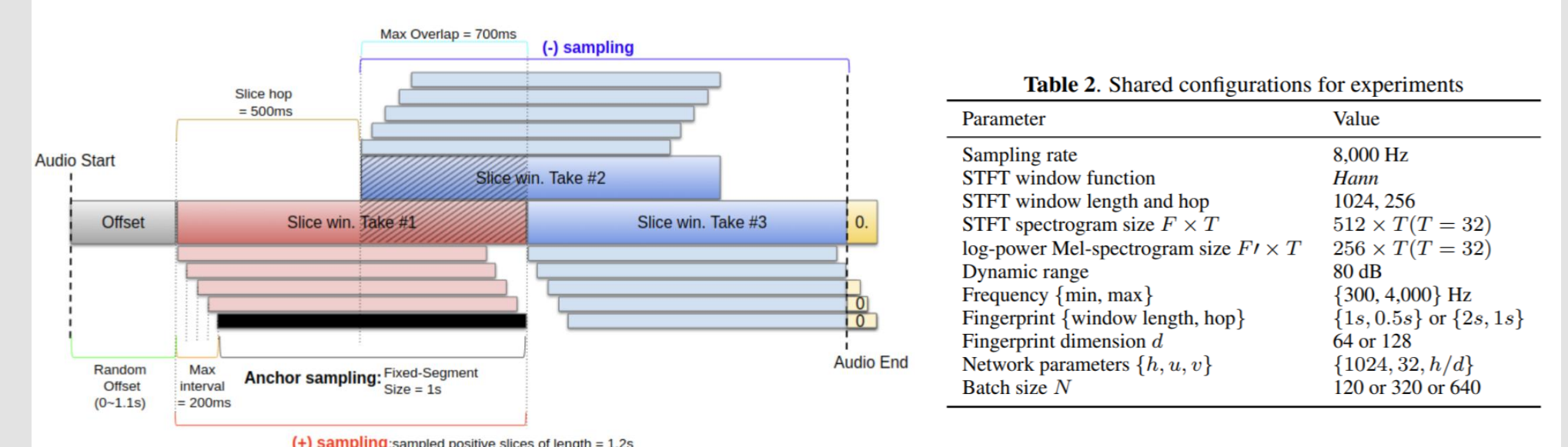
- The proposed contrastive learning framework could simulate the search task by explicitly sampling original-replica (clean-augmented) pairs as positive pairs.
- In the segment-level search task our model performed better than the model with triplet embeddings.
- Our model, using 10x less memory and shorter query length (< 3 s) than an existing rule-based algorithm, outperformed in song-level search.
- The superior performance of our model in the task is not due to any single design choice, but a combination of design choices.
- This study implies that the audio fingerprinting task would inherently have self-supervised learning potentials.

Implementation Notes

Audio augmentation chain

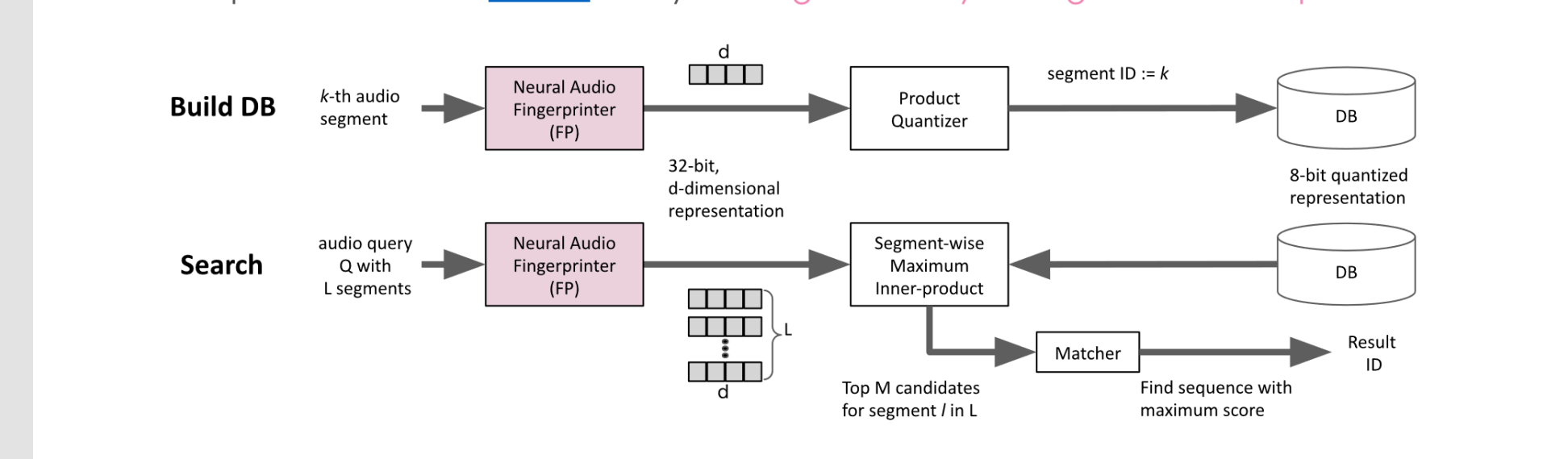


Time-offset modulation



Build DB & Search with sequence matcher

We build DB using only original audio sources. The resulting fingerprint vector is stored as a compressed representation using PQ (product quantizer) algorithm implemented with FAISS library. → **huge memory saving and search speed!**



References & Acknowledgements

- B. Gfeller et al. “Now playing: Continuous low-power music recognition,” in *NeurIPS 2017 Workshop on Machine Learning on the Phone and other Consumer Devices*, 2017.
- Chen, Ting, et al. “A simple framework for contrastive learning of visual representations.” *International conference on machine learning(ICML)*. 2020.
- <https://github.com/worldveil/dejavu>

We would like to thank the TensorFlow Research Cloud (TFRC) program that gave us access to Google Cloud TPUs.