

# Simultaneously Learning Robust Audio Embeddings and Balanced Hash Codes for Query-by-Example

Anup Singh<sup>1,2</sup>, Kris Demuyne<sup>1</sup>, Vipul Arora<sup>2</sup>

<sup>1</sup>Ghent University, <sup>2</sup>Indian Institute of Technology-Kanpur



## Introduction

### What is Audio Fingerprinting?

Audio fingerprinting generates compact summary of an audio signal.

### What are the applications?

Music identification, second-screen apps, broadcast monitoring, etc.

### What did existing methods do?

Conventional methods rely on handcrafted audio features to design audio fingerprints. Recent methods deploy deep learning to generate compact audio fingerprints.

### What makes it challenging?

- ✗ Robustness
- ✗ Efficient audio indexing
- ✗ Comprehensive audio search
- ✗ Practically deployable

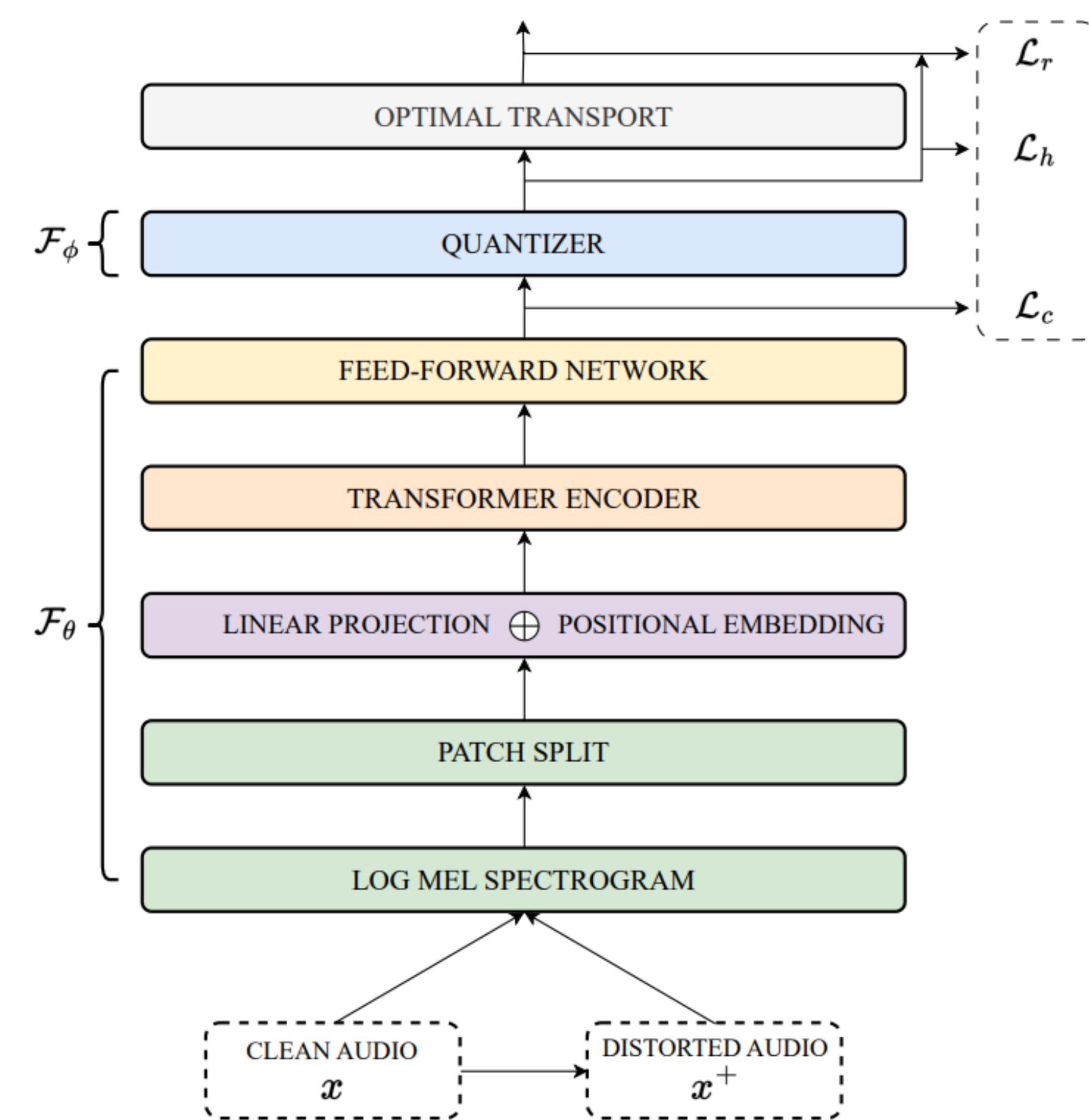
### What do we do?

- ✓ Discriminative embeddings using Transformer encoder
- ✓ Contrastive learning framework to achieve robustness
- ✓ Efficient indexing by learning balanced hash codes
- ✓ Adopted Optimal Transport to achieve balanced clustering
- ✓ Fine-grained audio search

### Research Questions:

- How to learn embeddings and hash codes to make audio indexing and retrieval efficient?
- How to achieve balanced clustering?

## Method



- Balanced clustering objective is combinatorial optimization problem - **Difficult to solve!**
- We adopt the optimal transport framework to obtain approximate solution.
- We also propose simple yet efficient subsequent search strategy to precisely locate the query timestamp.

### Training Objectives

$$\mathcal{L}_c = -\log \frac{e^{(\mathcal{F}_\theta(x) \cdot \mathcal{F}_\theta(x^+)) / \tau}}{e^{(\mathcal{F}_\theta(x) \cdot \mathcal{F}_\theta(x^+)) / \tau} + \sum_{x^-} e^{(\mathcal{F}_\theta(x) \cdot \mathcal{F}_\theta(x^-)) / \tau}}$$

$$\mathcal{L}_h = -\log \frac{e^{(q(x) \cdot q(x^+)) / \tau}}{e^{(q(x) \cdot q(x^+)) / \tau} + \sum_{x^-} e^{(q(x) \cdot q(x^-)) / \tau}}$$

$$\mathcal{L}_r = -[\text{softmax}(\frac{s_{k^*}(x)}{\tau}) + \text{softmax}(\frac{s_{k^*}(x^+)}{\tau})]$$

### Balanced Clustering Objective

$$\max_A \mathbb{E}_x \left[ \sum_{k=1}^{2^K} A_k(x) s_k(x) \right] \quad \text{subject to}$$

$$A_k(x) \in \{0, 1\}, \sum_{k=1}^{2^K} A_k(x) = 1 \quad \forall x, \mathbb{E}_x[A_k(x)] = \frac{1}{2^K} \quad \forall k$$

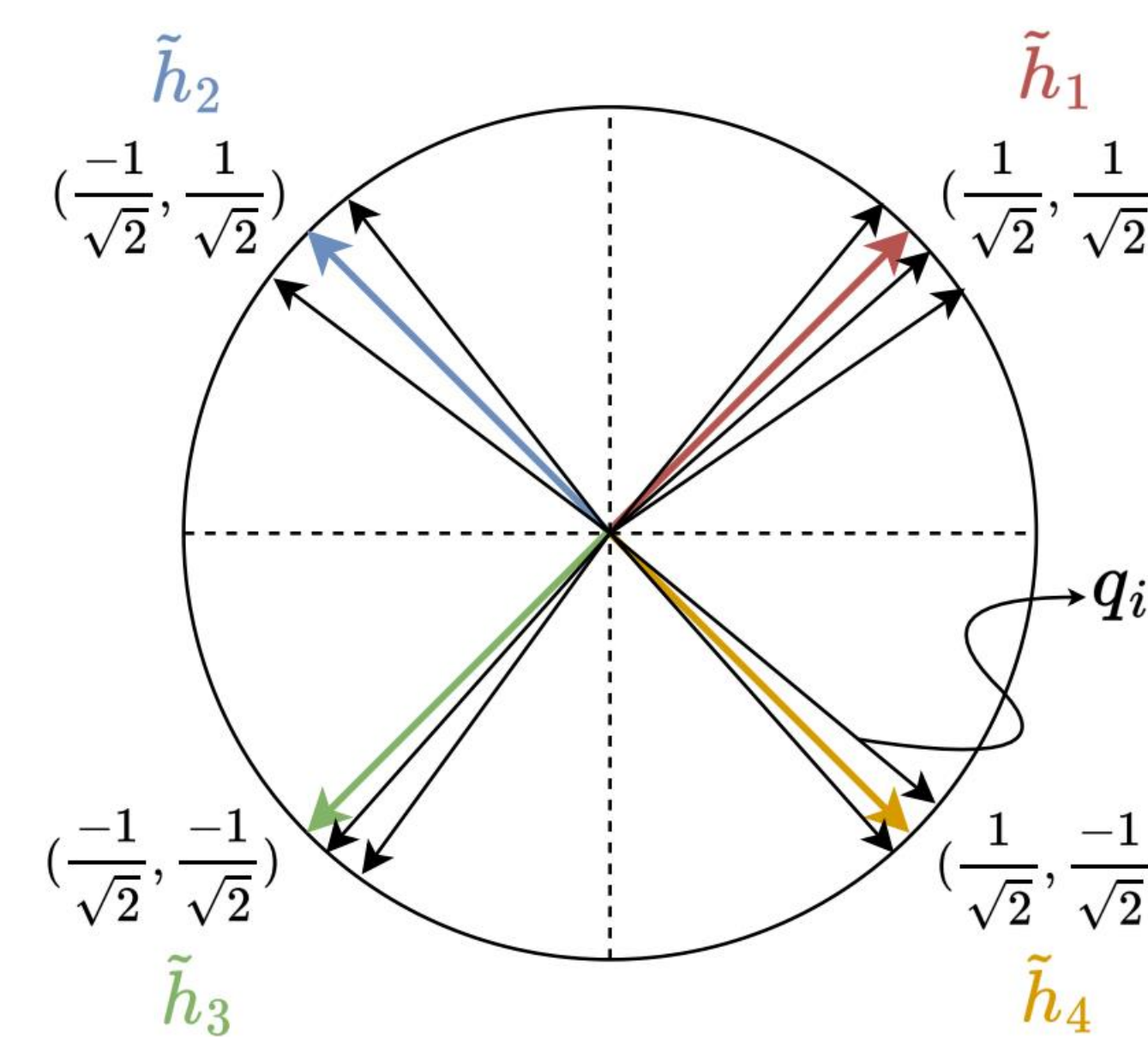


Illustration of balanced clustering in two-dimensional space

## Results

Length		Noise					Noise + Reverb					Reverb					speedup
		0	5	10	15	20	0	5	10	15	20	0.2	0.4	0.5	0.7	0.8	
1	NAFP + LSH	50.7	69.7	73.7	76.0	76.9	20.8	43.9	55.5	58.5	58.9	62.5	61.5	58.9	49.3	45.0	1x
	FE + LSH	<b>66.6</b>	<b>82.6</b>	87.6	90.0	90.6	<b>44.8</b>	<b>62.6</b>	73.8	<b>79.4</b>	<b>82.0</b>	83.8	83.8	78.4	<b>69.4</b>	64.4	1x
	FE+HT(Ours)	64.9	81.7	<b>87.9</b>	<b>90.1</b>	<b>90.9</b>	41.9	62.0	<b>74.6</b>	79.2	81.2	<b>84.6</b>	<b>82.3</b>	<b>80.0</b>	68.7	<b>64.5</b>	<b>2.4x</b>
2	NAFP + LSH	71.0	83.4	85.5	87.6	87.5	37.8	65.6	73.5	75.2	76.7	78.6	76.4	75.7	68.0	59.5	1x
	FE + LSH	<b>80.4</b>	88.2	91.6	93.2	94.8	<b>63.4</b>	78.2	84.6	86.0	85.4	90.4	86.4	85.0	74.8	67.8	1x
	FE+HT(Ours)	80.2	<b>89.4</b>	<b>93.1</b>	<b>94.1</b>	<b>94.9</b>	61.4	<b>80.8</b>	<b>85.0</b>	<b>87.5</b>	<b>88.3</b>	<b>92.4</b>	<b>87.7</b>	<b>86.4</b>	<b>77.3</b>	<b>70.7</b>	<b>2.4x</b>
3	NAFP + LSH	77.7	84.8	88.9	89.2	89.1	50.1	72.6	80.0	79.5	80.1	83.6	79.9	77.9	70.2	63.8	1x
	FE + LSH	83.2	88.4	92.6	94.4	95.2	<b>71.6</b>	82.6	85.6	86.2	87.4	91.8	87.8	85.6	74.4	68.4	1x
	FE+HT(Ours)	<b>84.7</b>	<b>90.8</b>	<b>95.5</b>	<b>96.3</b>	<b>97.1</b>	70.7	<b>84.1</b>	<b>88.6</b>	<b>89.0</b>	<b>90.1</b>	<b>93.8</b>	<b>90.3</b>	<b>87.9</b>	<b>77.9</b>	<b>71.6</b>	<b>2.3x</b>
5	NAFP + LSH	82.6	89.2	90.2	90.5	91.2	60.2	79.1	83.4	82.8	83.1	85.6	83.4	79.3	74.1	65.7	1x
	FE + LSH	85.6	90.0	92.8	94.2	95.8	80.0	87.0	87.1	87.6	87.2	93.8	88.8	86.6	75.0	69.0	1x
	FE+HT(Ours)	<b>88.0</b>	<b>93.4</b>	<b>95.3</b>	<b>96.2</b>	<b>97.4</b>	<b>80.6</b>	<b>88.6</b>	<b>90.8</b>	<b>91.3</b>	<b>91.5</b>	<b>96.4</b>	<b>91.1</b>	<b>88.9</b>	<b>78.7</b>	<b>71.2</b>	<b>2.4x</b>

## Experiments

- Database: Free Music Archive (FMA)
- Evaluation: Identify as a match if the located timestamp is within +/-50 ms.
- Baselines: NAFP<sup>1</sup> and Audfprint<sup>2</sup>
- We compare our indexing performance gain against the Locality Sensitive Hashing (LSH<sup>3</sup>)

## Findings

- Our **fingerprints are robust and discriminative** and thus achieve superior performance.
- Our system delivers **reliable performance with short query lengths** at high distortion levels.
- We achieve **more than two-folds search speedup** than LSH due to balanced hash codes.
- We achieve **higher efficiency as well efficacy** compared to baselines.

## References

1. S. Chang, D. Lee, J. Park, H. Lim, K. Lee, K. Ko, and Y. Han, "Neural audio fingerprint for high-specific audio retrieval based on contrastive learning," in ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2021, pp. 3025–3029.
2. <https://github.com/dpwe/audfprint>
3. Mayur Datar, Nicole Immerlica, Piotr Indyk, and Vahab S. Mirrokni, "Locality-sensitive hashing scheme based on p-stable distributions," in Proceedings of the twentieth annual symposium on Computational geometry, 2004, pp. 253–262.

Paper Link



Contact: anup.singh@ugent.be