

# FINGERPRINTING

## AUDIO AND VIDEO DETECTION SYSTEMS

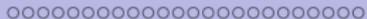
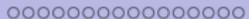
Speaker: Matteo Pozzetti

Supervisor: Laura Peer

01 April 2015

# OVERVIEW

- 1 INTRODUCTION
- 2 AUDIO FINGERPRINTING
- 3 VIDEO FINGERPRINTING
- 4 CONCLUSION



# Introduction



# FINGERPRINTING

Content based fingerprinting:

- Algorithm that maps a multimedia object into a compact digital summary (set of hashes)



# FINGERPRINTING

Content based fingerprinting:

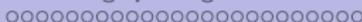
- Algorithm that maps a multimedia object into a compact digital summary (set of hashes)
- Used to retrieve such an object in a database



# FINGERPRINTING

## Content based fingerprinting:

- Algorithm that maps a multimedia object into a compact digital summary (set of hashes)
- Used to retrieve such an object in a database
- Applications:



# FINGERPRINTING

## Content based fingerprinting:

- Algorithm that maps a multimedia object into a compact digital summary (set of hashes)
- Used to retrieve such an object in a database
- Applications:
  - Search engines (Shazam)



# FINGERPRINTING

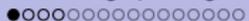
## Content based fingerprinting:

- Algorithm that maps a multimedia object into a compact digital summary (set of hashes)
- Used to retrieve such an object in a database
- Applications:
  - Search engines (Shazam)
  - Copyright protection (YouTube)



# TABLE OF CONTENTS

- 1 INTRODUCTION
- 2 AUDIO FINGERPRINTING**
  - Introduction
  - Hashing
  - Searching and Scoring
  - Results
- 3 VIDEO FINGERPRINTING
- 4 CONCLUSION



# Introduction





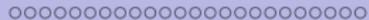
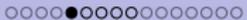
# GOALS

- The length of the recorded sample should not be more than 15 seconds long









# Hashing

# TRANSFORMATIONS TO THE FREQUENCY DOMAIN

## FOURIER TRANSFORM

Fourier transform: maps the signal  $x(t) : \mathbb{R} \rightarrow \mathbb{C}$  into

$$X(f) = \int_{-\infty}^{\infty} x(t) e^{-j2\pi ft} dt$$

which decomposes the function into its frequencies.



# SHORT-TIME FOURIER TRANSFORM

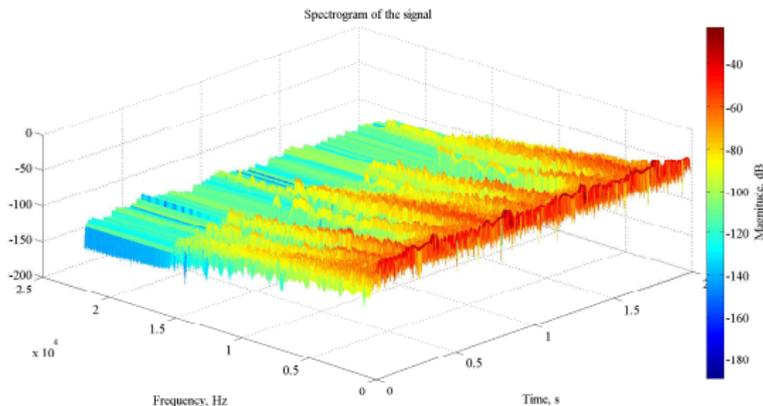
STFT is a Fourier-related transform used to determine the sinusoidal frequency and phase content of local sections of a signal as it changes over time.

$$x(t) \xrightarrow{STFT} X(\tau, f) = \int_{-\infty}^{\infty} x(t)w(t - \tau)e^{-j2\pi ft} dt$$

# SHORT-TIME FOURIER TRANSFORM

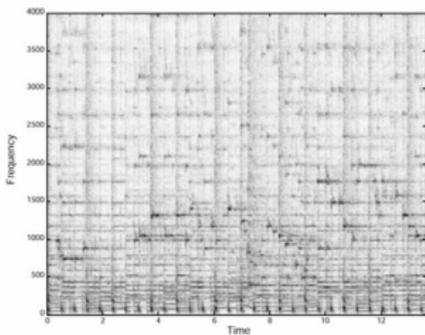
STFT is a Fourier-related transform used to determine the sinusoidal frequency and phase content of local sections of a signal as it changes over time.

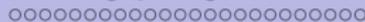
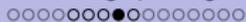
$$x(t) \xrightarrow{STFT} X(\tau, f) = \int_{-\infty}^{\infty} x(t)w(t - \tau)e^{-j2\pi ft} dt$$



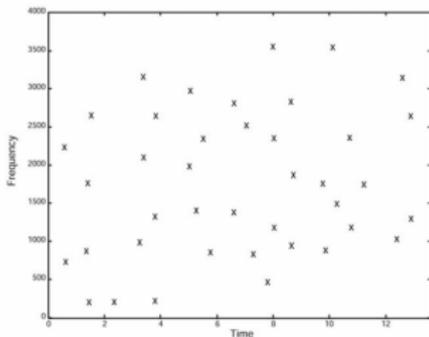
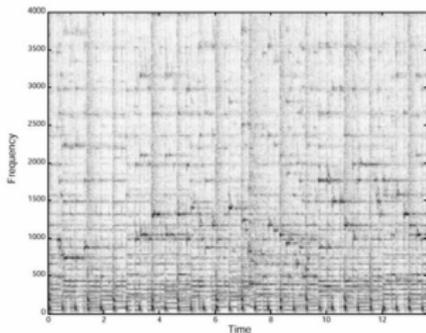


# CONSTELLATION AND HASHING

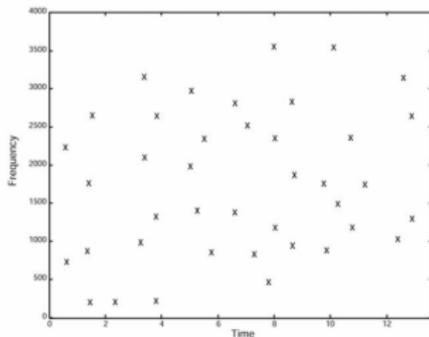
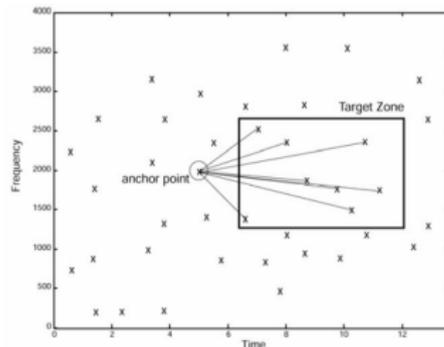
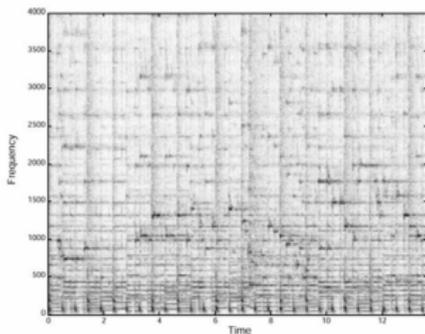


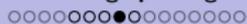


# CONSTELLATION AND HASHING

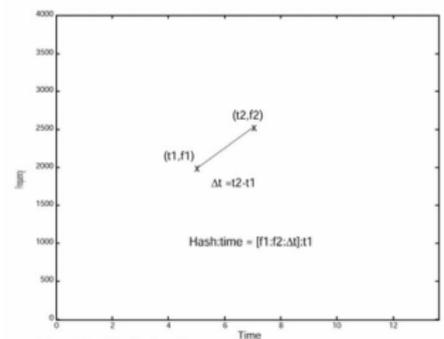
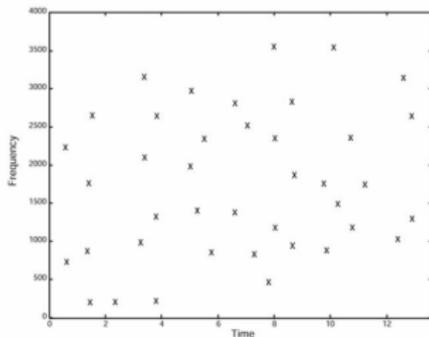
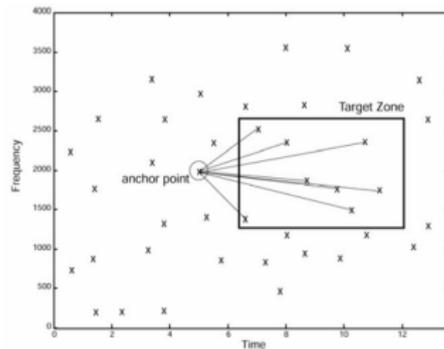
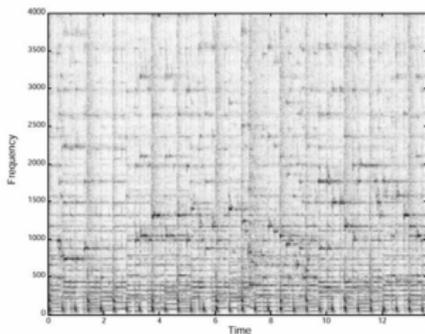


# CONSTELLATION AND HASHING





# CONSTELLATION AND HASHING



# HASHING

Each anchor point is sequentially paired with points within its target zone.

# HASHING

Each anchor point is sequentially paired with points within its target zone.

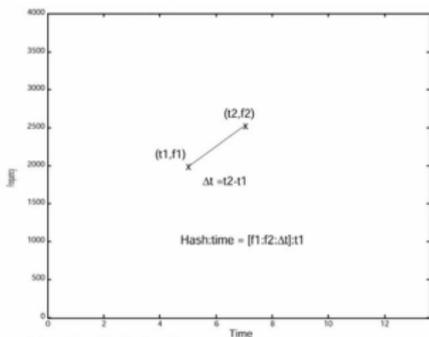
$$(f_1, f_2, \Delta t) \rightarrow \text{hash}$$

30 bits hash:

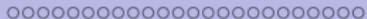
10 BITS frequency anchor

10 BITS frequency target

10 BITS time difference







## Searching and Scoring

# SEARCHING AND SCORING

We take all [matching hashes](#) of a given song and we plot their time offset

# SEARCHING AND SCORING

We take all **matching hashes** of a given song and we plot their time offset

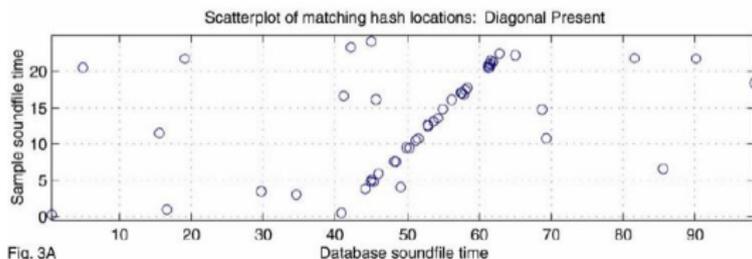
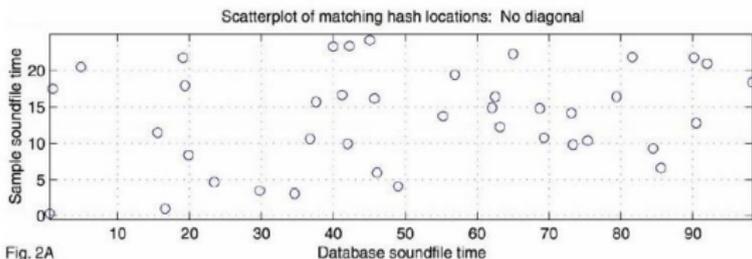
DB song	Query song
$h:t$	$h:\tau$



# SEARCHING AND SCORING

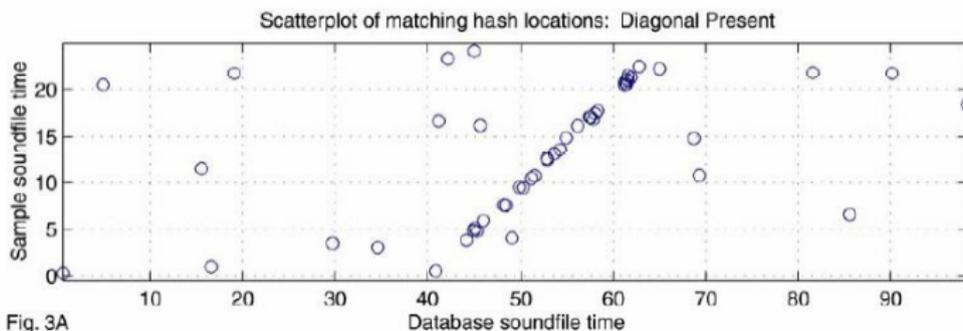
We take all **matching hashes** of a given song and we plot their time offset

DB song	Query song
$h:t$	$h:\mathcal{T}$



# SCORING

- Matching features:  $t_k = \tau_k + \text{constant}$



# SCORING

- Matching features:  $t_k = \tau_k + \text{constant}$
- For each  $(t_k, \tau_k)$  in the scatter-plot, we calculate:

$$\Delta_k = t_k - \tau_k$$

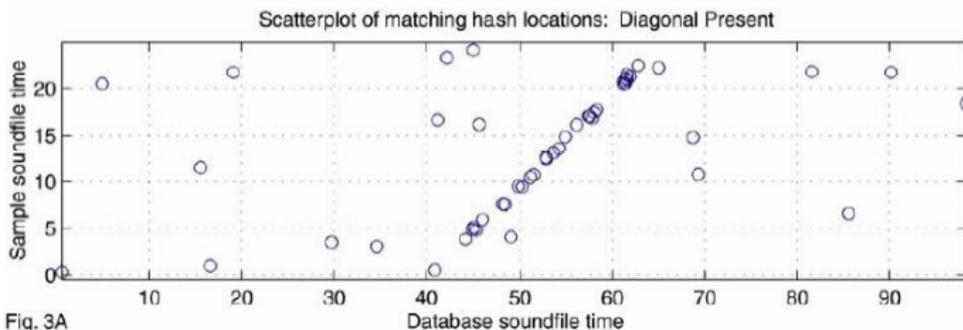


Fig. 3A

# SCORING

- Matching features:  $t_k = \tau_k + \text{constant}$
- For each  $(t_k, \tau_k)$  in the scatter-plot, we calculate:

$$\Delta_k = t_k - \tau_k$$

- Then we plot a histogram of these  $\Delta_k$

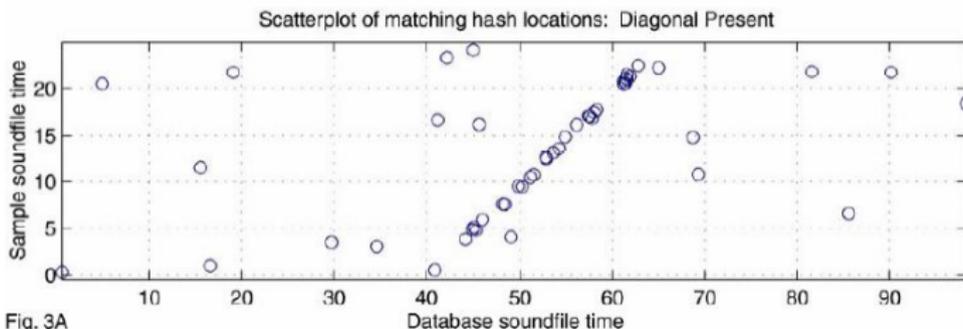
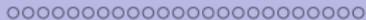
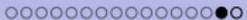


Fig. 3A



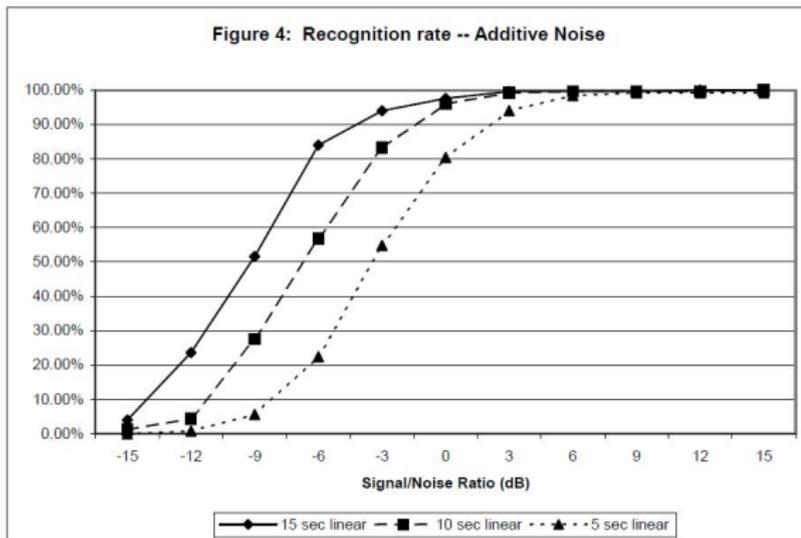




# Results

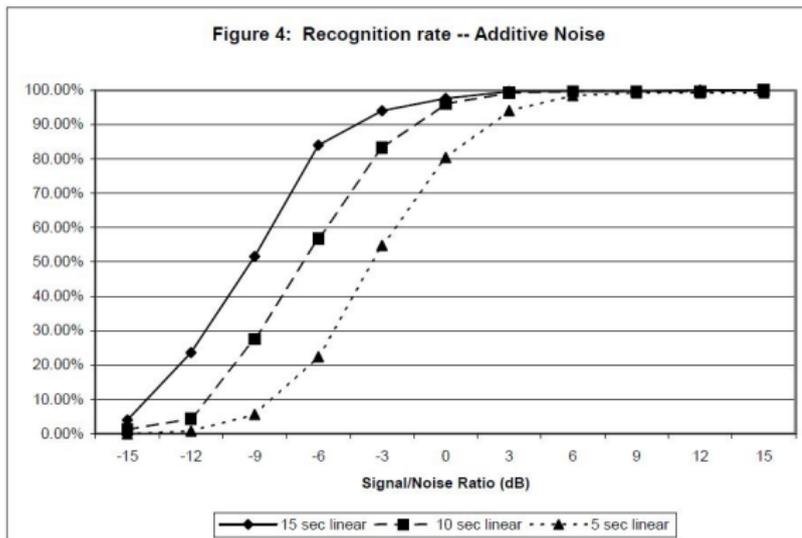
# RESULTS

- The algorithm performs well with significant levels of noise



# RESULTS

- The algorithm performs well with significant levels of noise



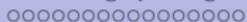
- For a database of 20 about thousand tracks the search time is in the order of 5 – 500 milliseconds.

# TABLE OF CONTENTS

- 1 INTRODUCTION
- 2 AUDIO FINGERPRINTING
- 3 VIDEO FINGERPRINTING
  - Introduction
  - Preprocessing
  - Discrete cosine transform
  - 3D-DCT
  - Fingerprinting using TIRIs
  - Searching
  - Results
- 4 CONCLUSION



# Introduction



# INTRODUCTION

## COPY DETECTION SYSTEM

The purpose of this system is to detect copyrights violations.

# INTRODUCTION

## COPY DETECTION SYSTEM

The purpose of this system is to detect copyrights violations.

A fingerprint should be

- Robust to content-preserving distortions



# INTRODUCTION

## COPY DETECTION SYSTEM

The purpose of this system is to detect copyrights violations.

A fingerprint should be

- Robust to content-preserving distortions
- Discriminant



# INTRODUCTION

## COPY DETECTION SYSTEM

The purpose of this system is to detect copyrights violations.

A fingerprint should be

- Robust to content-preserving distortions
- Discriminant
- Easy to compute



# INTRODUCTION

## COPY DETECTION SYSTEM

The purpose of this system is to detect copyrights violations.

A fingerprint should be

- Robust to content-preserving distortions
- Discriminant
- Easy to compute
- Compact



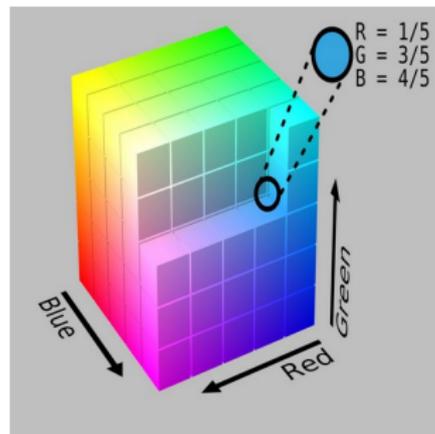




# IMAGE AND VIDEO REPRESENTATIONS

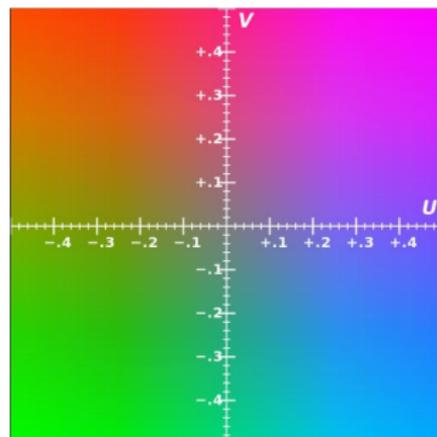
## • RGB

- Colors represented as a combination of red, green and blue (3D-space)
- Each component is an integer value in  $[0, 255]$



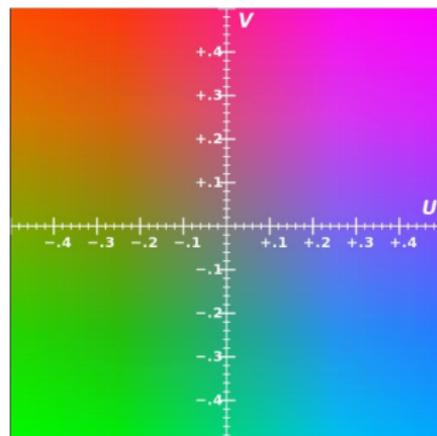
# IMAGE AND VIDEO REPRESENTATIONS

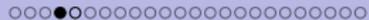
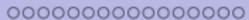
- RGB
  - Colors represented as a combination of red, green and blue (3D-space)
  - Each component is an integer value in  $[0, 255]$
- YUV



# IMAGE AND VIDEO REPRESENTATIONS

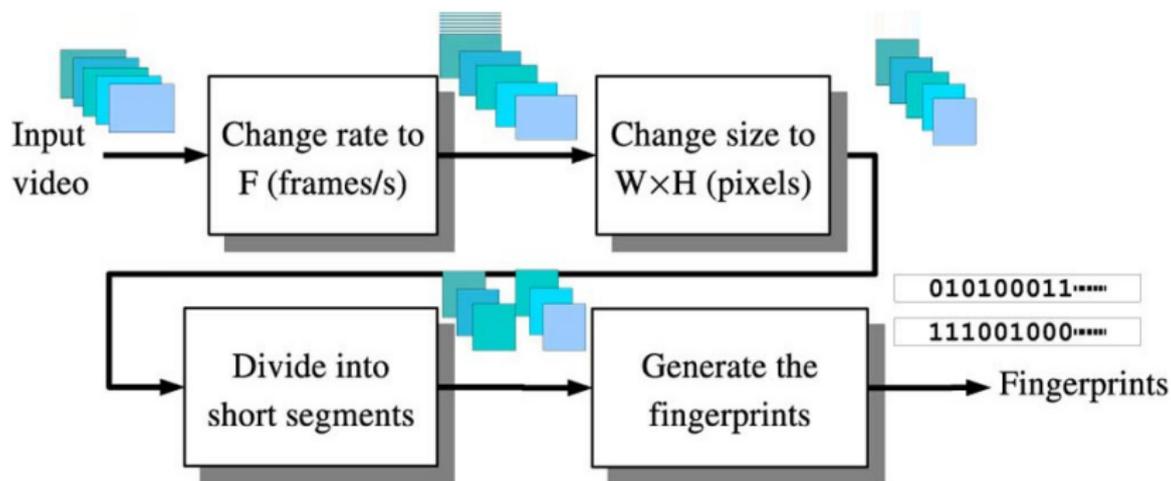
- RGB
  - Colors represented as a combination of red, green and blue (3D-space)
  - Each component is an integer value in  $[0, 255]$
- YUV
  - Y is the luminance component, U and V are the chrominance components.

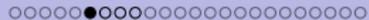
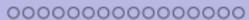




# Preprocessing

# PREPROCESSING





## Discrete cosine transform

# DISCRETE COSINE TRANSFORM

- The DCT is similar to the discrete Fourier transform: it transforms a signal or image from the spatial domain to the frequency domain.

# DISCRETE COSINE TRANSFORM

- The DCT is similar to the discrete Fourier transform: it transforms a signal or image from the spatial domain to the frequency domain.
- It is a widely used tool in image and video processing and compression!

# DISCRETE COSINE TRANSFORM

- The DCT is similar to the discrete Fourier transform: it transforms a signal or image from the spatial domain to the frequency domain.
- It is a widely used tool in image and video processing and compression!
- It is possible to extract some relevant information from the frequency representation

# DISCRETE COSINE TRANSFORM

## FORMULA

The two-dimensional DCT of an  $M$ -by- $N$  matrix  $A$  is defined as follow:

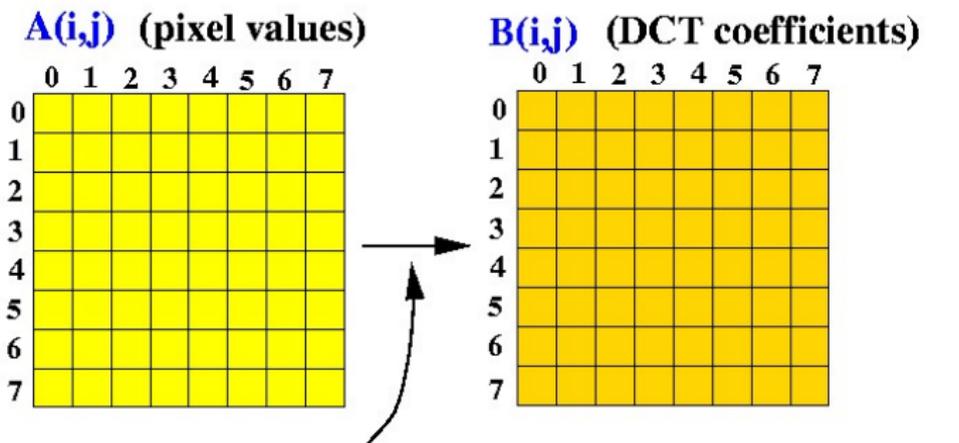
$$B_{pq} = \alpha_p \alpha_q \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} A_{mn} \cos \frac{\pi(2m+1)p}{2M} \cos \frac{\pi(2n+1)q}{2N}$$

$$\alpha_p = \begin{cases} \frac{1}{\sqrt{M}}, & p = 0 \\ \sqrt{\frac{2}{M}}, & 1 \leq p \leq M-1 \end{cases} \quad \alpha_q = \begin{cases} \frac{1}{\sqrt{N}}, & q = 0 \\ \sqrt{\frac{2}{N}}, & 1 \leq q \leq N-1 \end{cases}$$

where  $0 \leq p \leq M-1$  and  $0 \leq q \leq N-1$  are the matrix indexes.

## DISCRETE COSINE TRANSFORM

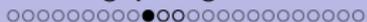
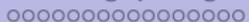
## EXAMPLE



$$B(k_1, k_2) = \frac{1}{4} C(k_1) C(k_2) \sum_{i=0}^7 \sum_{j=0}^7 A(i, j) \cos \left[ \frac{\pi \times k_1}{2 \times 8} (2i + 1) \right] \cos \left[ \frac{\pi \times k_2}{2 \times 8} (2j + 1) \right]$$

$$C(k) = \frac{1}{\sqrt{2}} \text{ if } k = 0$$

$$C(k) = 1, \text{ otherwise}$$



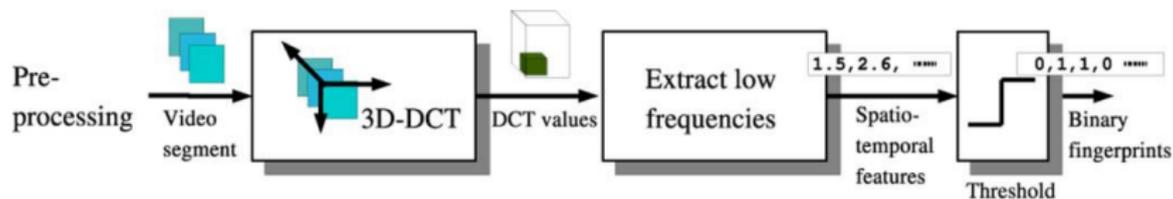
## 3D-DCT

# 3D DISCRETE COSINE TRANSFORM

3D-DCT is a DCT applied to a 3-dimensional matrix.

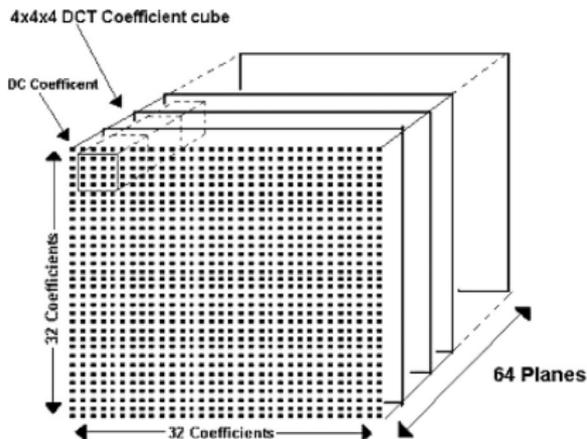
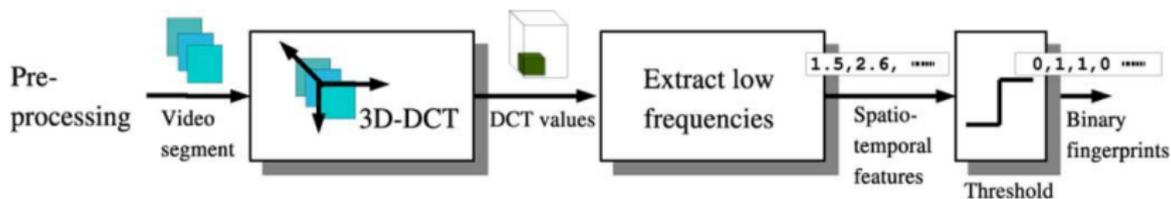


# FINGERPRINTING USING 3D-DCT

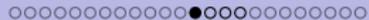




# FINGERPRINTING USING 3D-DCT



Binary fingerprints are derived by thresholding the low-frequency coefficients of the transform; the **threshold** is the **median** value of the selected coefficients.



# Fingerprinting using TIRIs

# TIRI

TEMPORALLY INFORMATIVE REPRESENTATIVE IMAGE

A TIRI contains spatial and temporal information of a short segment of a video sequence.

# TIRI

## TEMPORALLY INFORMATIVE REPRESENTATIVE IMAGE

A TIRI contains spatial and temporal information of a short segment of a video sequence.



$$I'_{m,n} = \sum_{k=1}^J w_k I_{m,n,k}$$

where  $I_{m,n,k}$  is the luminance value of the  $(m,n)$ th pixel in a set of  $J$  frames and  $w_k$  is a weighting factor.

# TIRI

## WEIGHTING FACTOR

Experiments were made with different  $w$ : constant, linear, exponential.



(a)



(b)



(c)



(d)



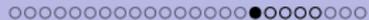
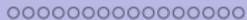
(e)



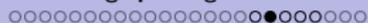
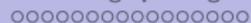
(f)

FIGURE : (d)  $w_k = 1$  (constant), (e)  $w_k = k$  (linear), (f)  $w_k = 1.2^k$  (exponential)





# Searching



# SEARCHING REQUIREMENTS

After the fingerprint is extracted, the database is searched for the closest fingerprint.

# SEARCHING REQUIREMENTS

After the fingerprint is extracted, the database is searched for the closest fingerprint.

## NEAREST NEIGHBOR PROBLEM

Fingerprints of two different copies of the same video are similar but not necessarily identical. We then search for the most similar fingerprint in the binary space.

# SEARCHING REQUIREMENTS

After the fingerprint is extracted, the database is searched for the closest fingerprint.

## NEAREST NEIGHBOR PROBLEM

Fingerprints of two different copies of the same video are similar but not necessarily identical. We then search for the most similar fingerprint in the binary space.

Since the database is very large, the algorithm must be fast. A *fast approximate algorithm* is preferred over an *exact* slow one.

# SEARCHING REQUIREMENTS

After the fingerprint is extracted, the database is searched for the closest fingerprint.

## NEAREST NEIGHBOR PROBLEM

Fingerprints of two different copies of the same video are similar but not necessarily identical. We then search for the most similar fingerprint in the binary space.

Since the database is very large, the algorithm must be fast. A *fast approximate algorithm* is preferred over an *exact* slow one.

## HAMMING DISTANCE

The process used to create binary fingerprints allow us to use the Hamming distance as a metric of similarity.

# INVERTED-FILE-BASED SIMILARITY SEARCH

$\underbrace{0110001}_{1^{st} \text{ word}}$ 
 $\underbrace{0010110}_{2^{nd} \text{ word}}$ 
 $\dots$ 
 $\underbrace{1110001}_{(n-1)^{th} \text{ word}}$ 
 $\underbrace{0010101}_{n^{th} \text{ word}}$

$m = \text{word length}$

$$n = \frac{L}{m}$$





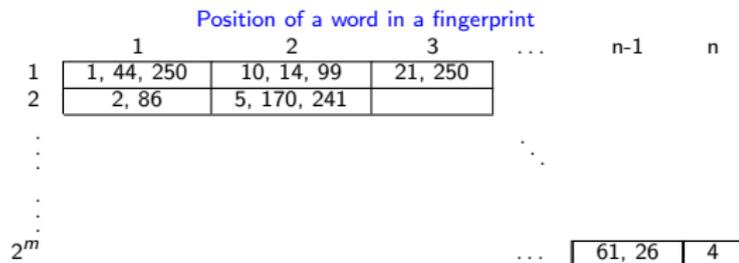


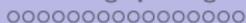


# INVERTED-FILE-BASED SIMILARITY SEARCH

To find a query fingerprint in a database:

- The fingerprint is divided into  $n$  words (of  $m$  bits)
- The query is compared with fingerprints that starts with the same word.
- If the Hamming distance is below a certain threshold, the match is found.
- Otherwise the procedure is repeated with the second word and so on.





# CLUSTER-BASED SIMILARITY SEARCH

- The database is clustered into  $K$  non-overlapping groups; each cluster has a centroid called **cluster head**.

# CLUSTER-BASED SIMILARITY SEARCH

- The database is clustered into  $K$  non-overlapping groups; each cluster has a centroid called **cluster head**.
- The cluster head closest to the query fingerprint is found

$\underbrace{0100001}_0$ 
 $\underbrace{0010110}_0$ 
 $\underbrace{1010111}_1$ 
 $\underbrace{0000001}_0$ 
 $\underbrace{1111101}_1$   
 $\text{dist}_0 = 2$     $\text{dist}_0 = 3$     $\text{dist}_1 = 2$     $\text{dist}_0 = 1$     $\text{dist}_1 = 1$

Cluster head = 00101

# CLUSTER-BASED SIMILARITY SEARCH

- The database is clustered into  $K$  non-overlapping groups; each cluster has a centroid called **cluster head**.
- The cluster head closest to the query fingerprint is found

$\underbrace{0100001}_0 \underbrace{0010110}_0 \underbrace{1010111}_1 \underbrace{0000001}_0 \underbrace{1111101}_1$   
 $\text{dist}_0 = 2 \quad \text{dist}_0 = 3 \quad \text{dist}_1 = 2 \quad \text{dist}_0 = 1 \quad \text{dist}_1 = 1$

Cluster head = 00101

- Every fingerprint belonging to this cluster are searched to find a match

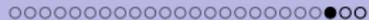
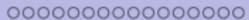
# CLUSTER-BASED SIMILARITY SEARCH

- The database is clustered into  $K$  non-overlapping groups; each cluster has a centroid called **cluster head**.
- The cluster head closest to the query fingerprint is found

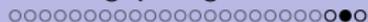
$\underbrace{0100001}_0$ 
 $\underbrace{0010110}_0$ 
 $\underbrace{1010111}_1$ 
 $\underbrace{0000001}_0$ 
 $\underbrace{1111101}_1$   
 $\text{dist}_0 = 2$     $\text{dist}_0 = 3$     $\text{dist}_1 = 2$     $\text{dist}_0 = 1$     $\text{dist}_1 = 1$

Cluster head = 00101

- Every fingerprint belonging to this cluster are searched to find a match
- If no match is found, then the cluster that is the second closest to the query is examined, then the third and so on



# Results



# RESULTS

- A database of 200 videos was created
- Videos were attacked (distorted) to create query videos

# RESULTS

- A database of 200 videos was created
- Videos were attacked (distorted) to create query videos

## METRICS

**RECALL** It's the true positive rate. It is a measure of robustness.

$$\text{Recall} = \frac{TP}{TP + FN}$$

# RESULTS

- A database of 200 videos was created
- Videos were attacked (distorted) to create query videos

## METRICS

**RECALL** It's the true positive rate. It is a measure of robustness.

$$\text{Recall} = \frac{TP}{TP + FN}$$

**PRECISION** Percentage of correct hits within all detected copies.

$$\text{Precision} = \frac{TP}{TP + FP}$$

# RESULTS

- A database of 200 videos was created
- Videos were attacked (distorted) to create query videos

## METRICS

**RECALL** It's the true positive rate. It is a measure of robustness.

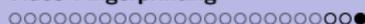
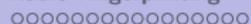
$$\text{Recall} = \frac{TP}{TP + FN}$$

**PRECISION** Percentage of correct hits within all detected copies.

$$\text{Precision} = \frac{TP}{TP + FP}$$

**F-SCORE** It is the harmonic mean of precision and recall:

$$F = 2 \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$



# RESULTS

Attack	F-score	
	3D-DCT	TIRI-DCT
Noise	1.00	0.99
Brightness	1.00	0.99
Contrast	1.00	0.99
Rotation	1.00	1.00
Time shift	0.91	0.98
Spatial shift	1.00	0.98
Frame loss	0.98	0.99
Average	0.99	0.99

# RESULTS

Attack	F-score	
	3D-DCT	TIRI-DCT
Noise	1.00	0.99
Brightness	1.00	0.99
Contrast	1.00	0.99
Rotation	1.00	1.00
Time shift	0.91	0.98
Spatial shift	1.00	0.98
Frame loss	0.98	0.99
Average	0.99	0.99

## SPEED

- TIRI-DCT is more than 3 times faster than 3D-DCT

## RESULTS

Attack	F-score	
	3D-DCT	TIRI-DCT
Noise	1.00	0.99
Brightness	1.00	0.99
Contrast	1.00	0.99
Rotation	1.00	1.00
Time shift	0.91	0.98
Spatial shift	1.00	0.98
Frame loss	0.98	0.99
Average	0.99	0.99

## SPEED

- TIRI-DCT is more than 3 times faster than 3D-DCT
- Approximate search algorithms are much faster than the exhaustive search when the error rates are low.

# RESULTS

Attack	F-score	
	3D-DCT	TIRI-DCT
Noise	1.00	0.99
Brightness	1.00	0.99
Contrast	1.00	0.99
Rotation	1.00	1.00
Time shift	0.91	0.98
Spatial shift	1.00	0.98
Frame loss	0.98	0.99
Average	0.99	0.99

## SPEED

- TIRI-DCT is more than 3 times faster than 3D-DCT
- Approximate search algorithms are much faster than the exhaustive search when the error rates are low. Cluster-based approach is faster.

# TABLE OF CONTENTS

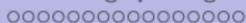
- 1 INTRODUCTION
- 2 AUDIO FINGERPRINTING
- 3 VIDEO FINGERPRINTING
- 4 CONCLUSION

# SHAZAM

- The algorithm is noise resistant and efficient. It can identify a short segment of music out of a database of over a million tracks. However the algorithm **is not meant to be resistant to general attacks** (e.g time compression, pitch alterations, ...).

# SHAZAM

- The algorithm is noise resistant and efficient. It can identify a short segment of music out of a database of over a million tracks. However the algorithm **is not meant to be resistant to general attacks** (e.g time compression, pitch alterations, ...).
- Shazam was founded in 1999 and launched its first identification service in 2002



# SHAZAM

- The algorithm is noise resistant and efficient. It can identify a short segment of music out of a database of over a million tracks. However the algorithm **is not meant to be resistant to general attacks** (e.g time compression, pitch alterations, ...).
- Shazam was founded in 1999 and launched its first identification service in 2002
- Shazam has now 400 million users in 200 countries and was used to identify 15 billion songs

# SHAZAM

- The algorithm is noise resistant and efficient. It can identify a short segment of music out of a database of over a million tracks. However the algorithm **is not meant to be resistant to general attacks** (e.g time compression, pitch alterations, ...).
- Shazam was founded in 1999 and launched its first identification service in 2002
- Shazam has now 400 million users in 200 countries and was used to identify 15 billion songs
- The value of the company is now half a billion dollars

# VIDEO FINGERPRINTING

- The proposed system is robust and it is able to find a match in a database of tens of millions of fingerprints in a few seconds

# VIDEO FINGERPRINTING

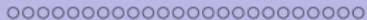
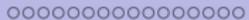
- The proposed system is robust and it is able to find a match in a database of tens of millions of fingerprints in a few seconds
- It yields a high average true positive rate of 97.6% and a low average false positive rate of 1%.

# YOUTUBE TESTING

Results of [Scott Smitelli](#) tests:

Modification Performed	Result
Song reversed	Pass
Pitch was lowered 6%	Pass
Pitch was lowered 5%	Fail
Pitch was raised 5%	Fail
Pitch was raised 6%	Pass
Time was expanded 6%	Pass
Time was expanded 5%	Fail
Time was compressed 5%	Fail
Time was compressed 6%	Pass
44% noise, 56% song	Fail
45% noise, 55% song	Pass

Modification Performed	Result
Volume was reduced 48 dB	Fail
Volume was reduced 18 dB	Fail
Volume was reduced 6 dB	Fail
Volume was increased 6 dB	Fail
Volume was increased 18 dB	Fail
Volume was increased 48 dB	Fail
Slow down 4%	Pass
Slow down 3%	Fail
Speed up 3%	Fail
Speed up 4%	Fail
Speed up 5%	Pass



# Questions?