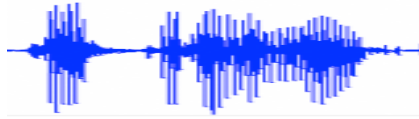


MASK: Robust Local Features for Audio Fingerprinting

Xavier Anguera, Antonio Garzón and
Tomasz Adamek
Telefonica Research

Outline

- What is audio fingerprinting
- MASK proposal
- Experiments
- Conclusions

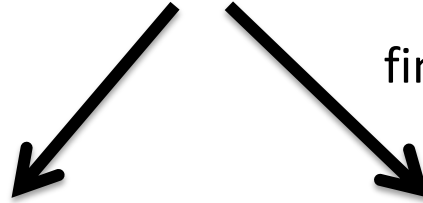


?



watermarking

fingerprinting

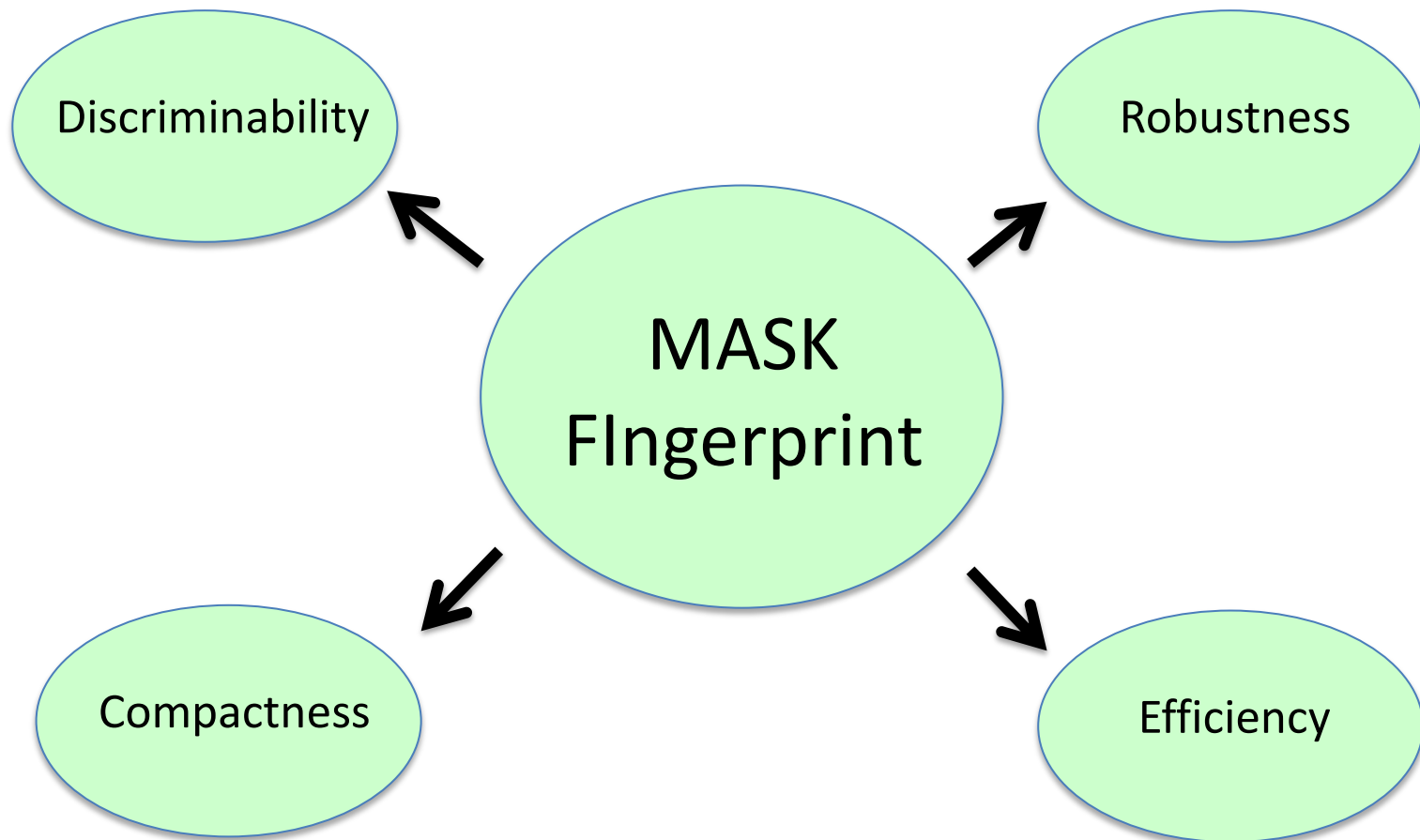




What makes a good audio fingerprint?



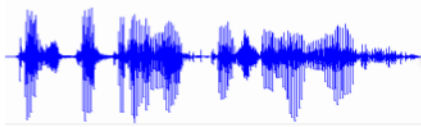
MusicBrainz



MASK == Masked Audio Spectral Keypoints

Considered prior art

- Avery Wang, “An industrial strength audio search algorithm,” in Proc. ISMIR, 2003.
- Jaap Haitsma and Antonius Kalker, “A highly robust audio fingerprinting system,” in Proc. ISMIR, 2002.
- Shumeet Baluja and Michele Covell, “Waveprint: Efficient wavelet-based audio fingerprinting”, Proc. Pattern Recognition 41 (2008)



Time-to-frequency
transformation



Salient spectral points
search



Local mask application

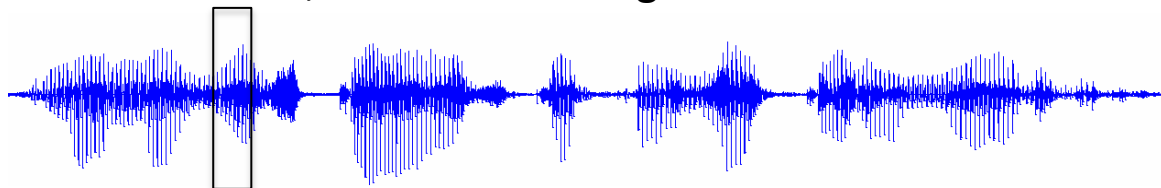


MASK fingerprint
encoding and storage

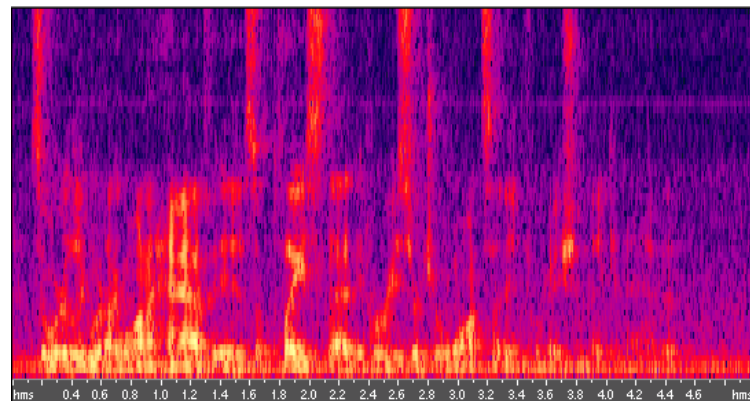
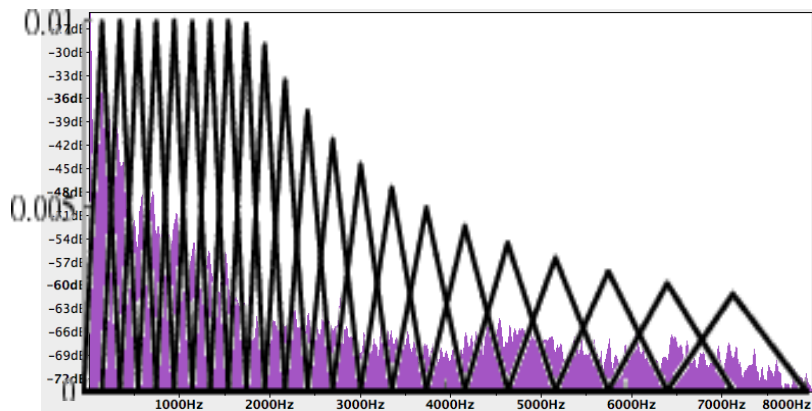


```
0110101000110110  
0001101110110111  
0001001001111010  
0000101111101010  
1111011000100110  
0010010100110101
```

10ms, 100ms Hamming window

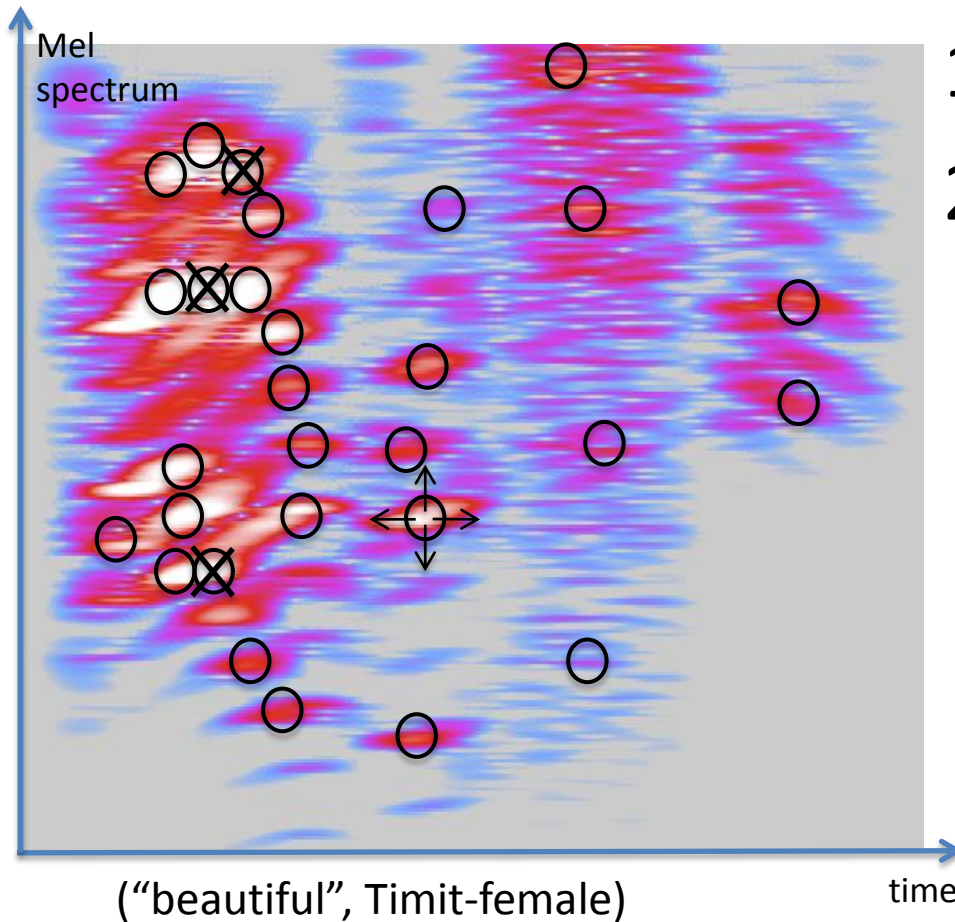


FFT 1024, bandwidth limited to 300-3KHz



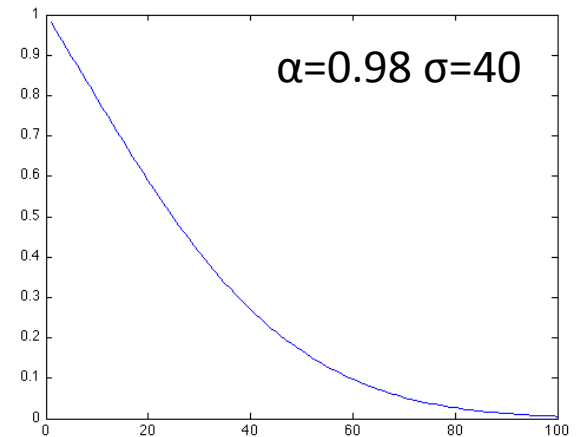
18 or 34 MEL-spectrum bands

Selection of salient spectral points

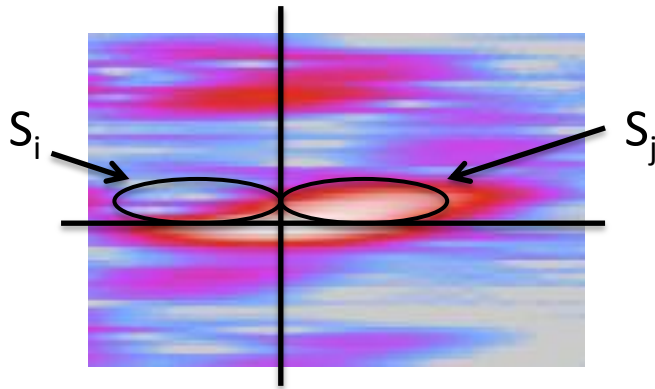


1. Detect all maximas
2. Trim to the desired number

$$Thr[n] = \alpha^{\Delta t} E[n - 1] \exp - \frac{(\Delta t)^2}{2 * \sigma^2}$$



Spectral masking around salient points

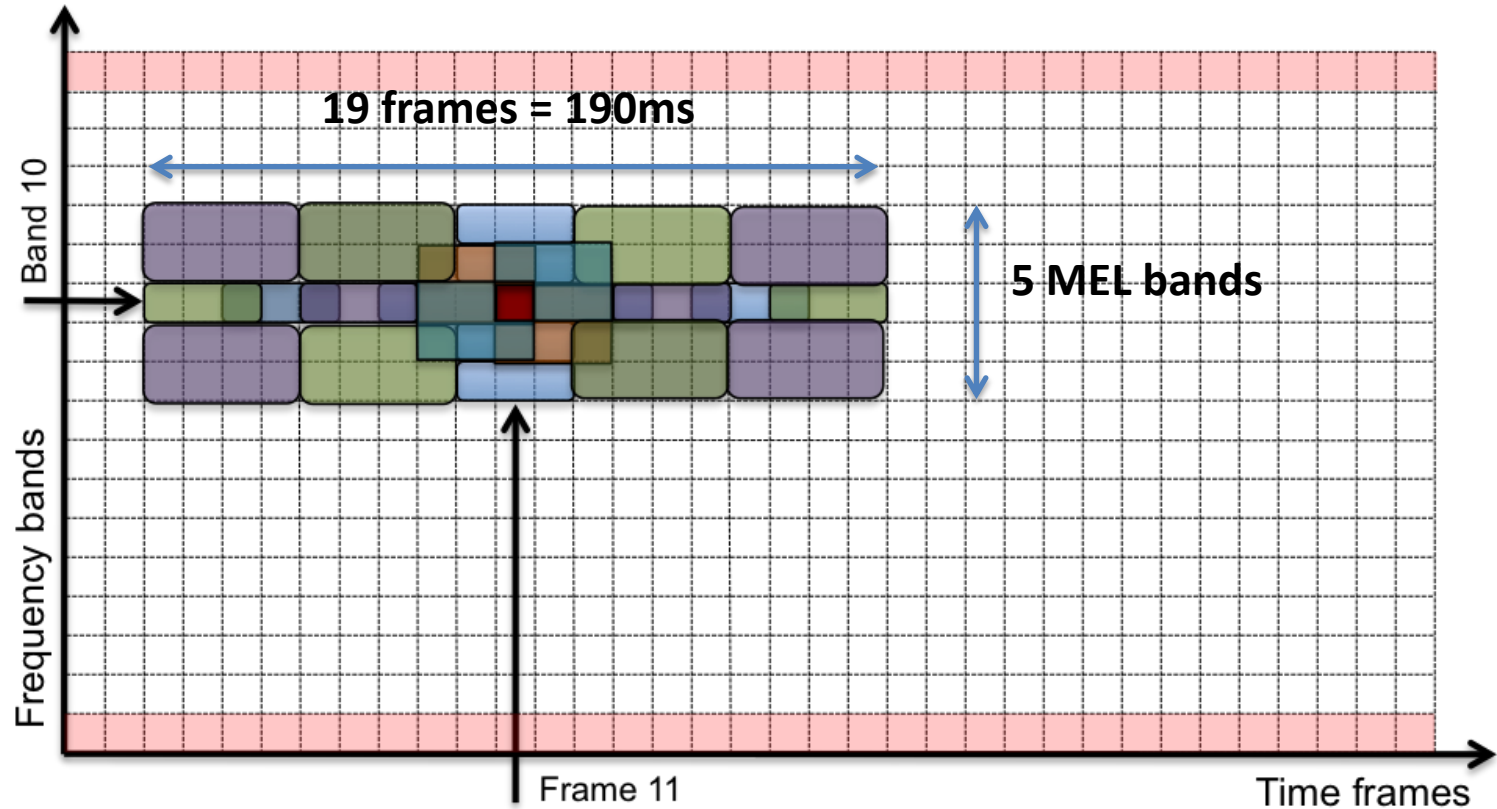


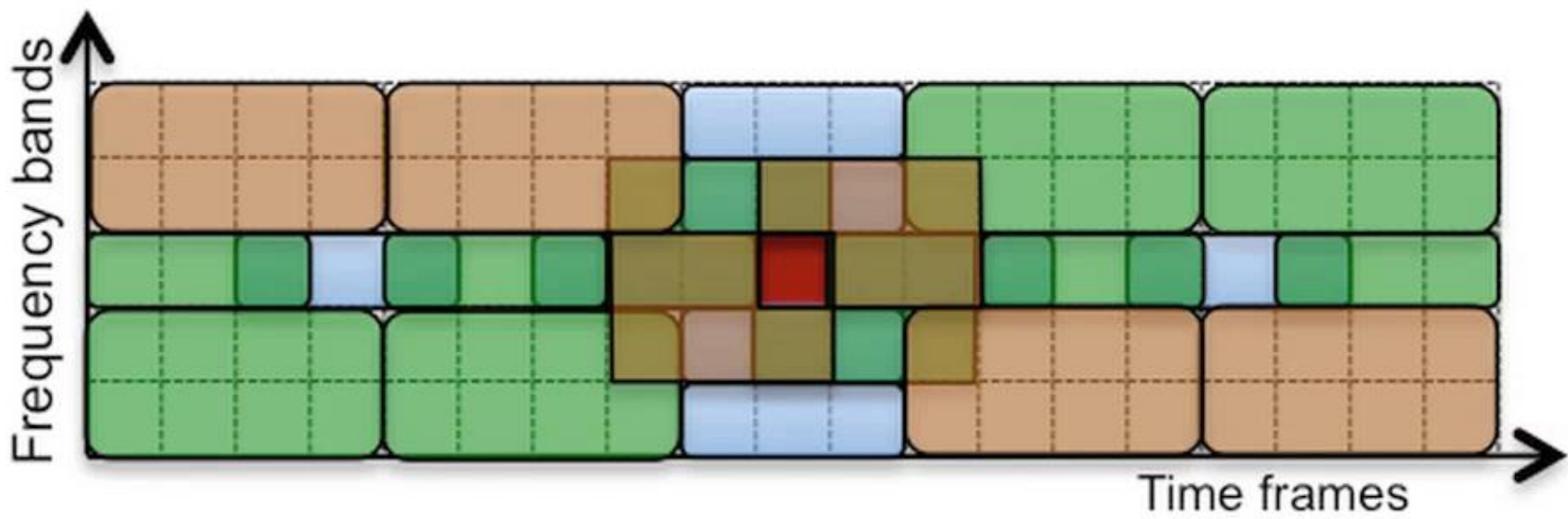
Spectral Regions

- Include one or several time-frequency values
- Overlaps are allowed
- The number of comparisons defines the size of the fingerprint
- Designed manually (for now)

$$\begin{cases} S_i > S_j \rightarrow b[n] = 1 \\ \textit{otherwise} \rightarrow b[n] = 0 \end{cases}$$

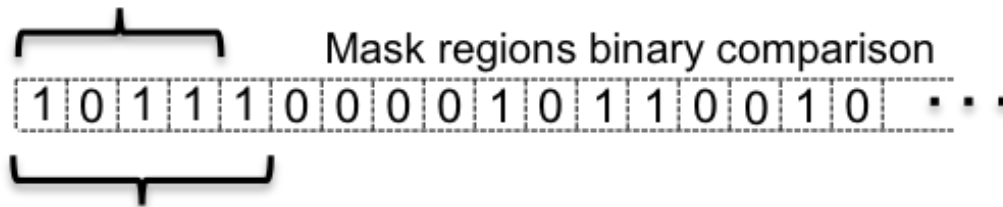
Current MASK regions





Fingerprint encoding

Peak location (16 bands)



Peak location (32 bands)

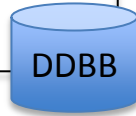
- 4-5 bits for the MEL band where the maxima is located
- 22 Bits obtained from spectral regions comparison

Indexing and retrieval

Reference inverted file index

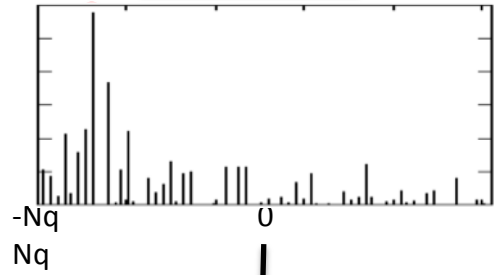
011...00	→	(movie1, 10); (movie6, 4); (movie9, 1); (movie7, 34)
001...01	→	(movie5, 35); (movie7, 80)
...	→	...
100...11	→	(movie9, 24); (movie3, 5); (movie8, 11)

MASK FP Content ID Time offset

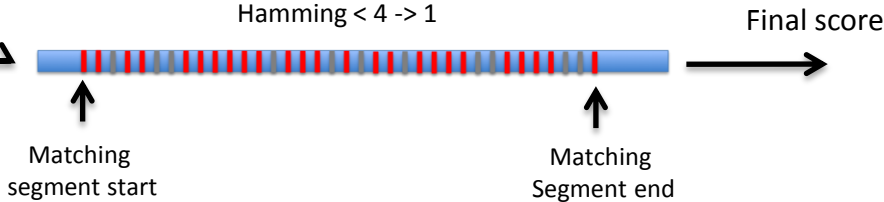


Exact matching

QUERY MASK
0->011...00
1->100...11
...
13->111...01
14->011...01
15->000...10



Hamming > 4 -> 0
Hamming < 4 -> 1



Experimental section



database

NIST-TRECVID 2010-2011 data for video-copy detection

- 400h reference videos
- 1400 audio queries per year (201 unique videos X 7)
- 7 audio transformations
 - original
 - MP3 compression
 - MP3 compression + multiband companding
 - Bandwidth limit (500-3K) + single-band companding
 - Mixed with speech
 - Mixed with speech + multiband companding
 - Mixed with speech + bandwidth limit + monoband companding

Metric & baseline

- normalized detection cost rate (NDCR) in balanced profile

$$NDCR_{BALANCED} = P_{MISS} + 200R_{FA}$$

- Compare results with a similar fingerprint to the Philips fingerprint

Jaap Haitsma and Antonius Kalker, "A highly robust audio fingerprinting system," in Proc. International Symposium on Music Information Retrieval (ISMIR), 2002.

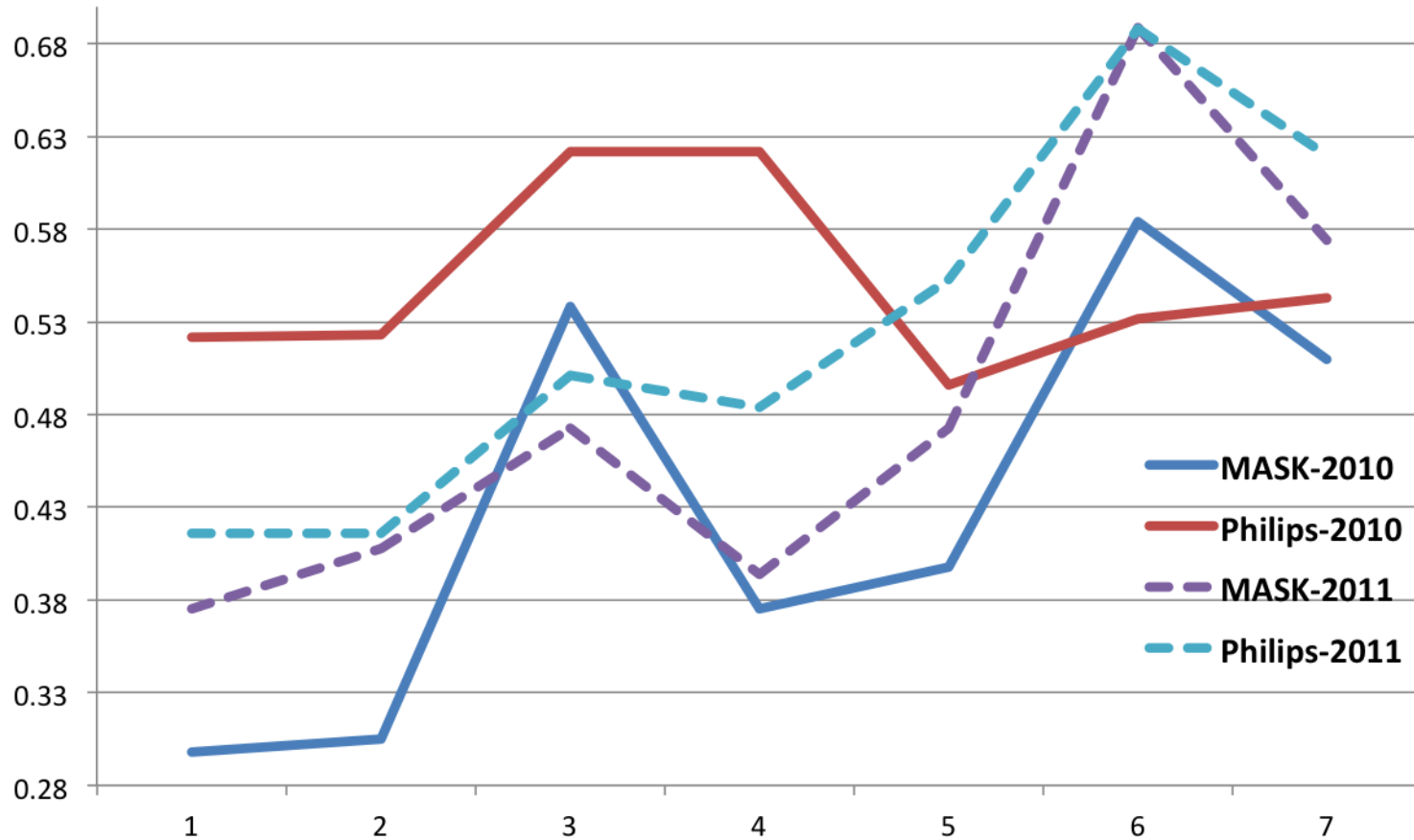
Table 2. Comparison of Minimum NDCR scores

system	# results	dataset	Min. NDCR	% improve.
Philips MASK	1	2010	0.55 0.43	— 21.8%
Philips MASK	20	2010	1.03 0.79	— 23.3%
Philips MASK	1	2011	0.53 0.48	— 9.4%
Philips MASK	20	2011	0.96 0.82	— 14.5%

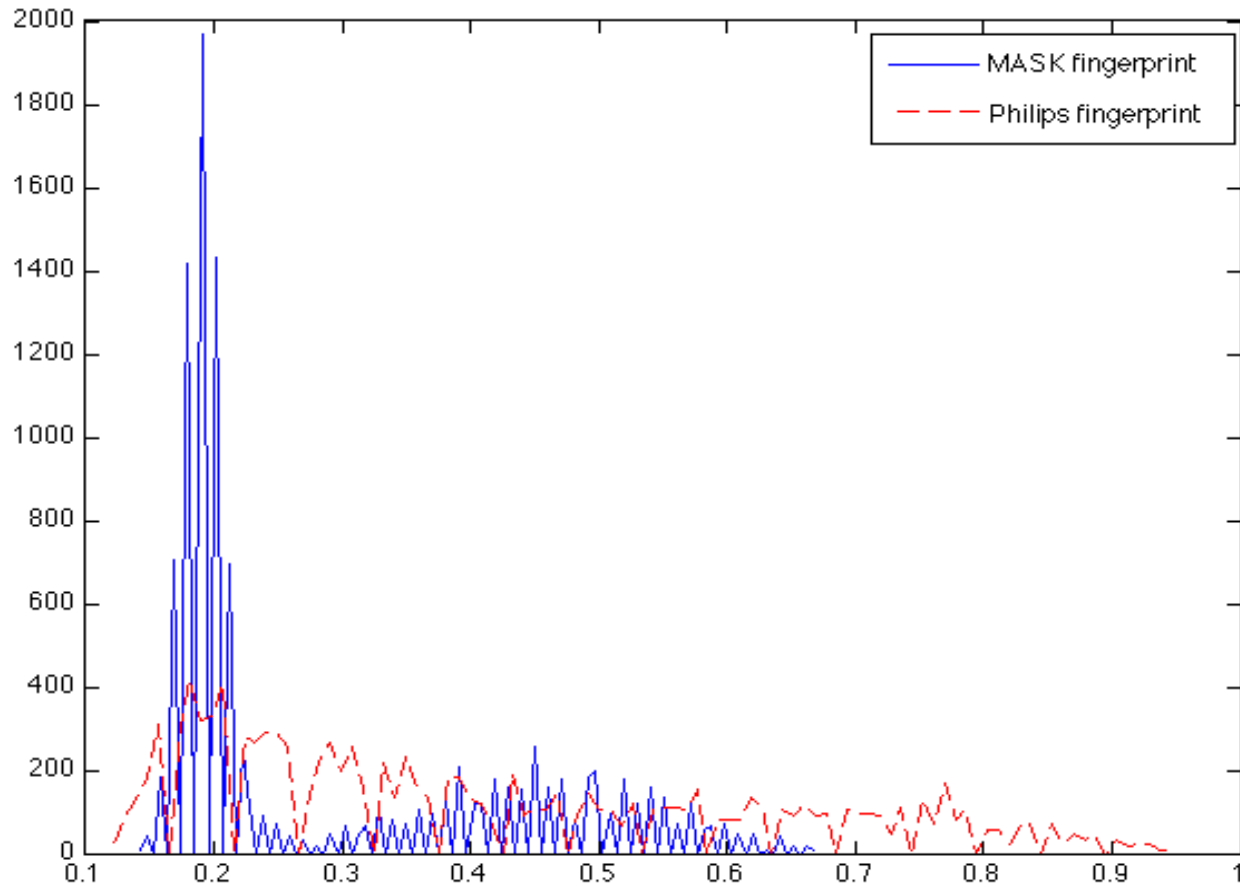
Table 3. Comparison of Actual NDCR scores

system	# results	dataset	Thr. std.	Act. NDCR	% improve.
Philips MASK	1	2010	0.11 0.03	0.60 0.44	– 26.6%
Philips MASK	20	2010	0.12 0.04	1.19 0.91	– 23.5%
Philips MASK	1	2011	0.08 0.03	0.57 0.50	– 12.2%
Philips MASK	20	2011	0.09 0.06	1.18 1.02	– 13.5%

Comparison per transformation



Scores histogram



Conclusions

- A novel binary fingerprint is proposed to improve on some shortcomings from well reputed prior art.
- We show that we can extract the FP and use it for VCD with excellent results