

FELIX STÜRMER
SKETCH-BASED IMAGE RETRIEVAL USING
CURVELETS

SKETCH-BASED IMAGE RETRIEVAL USING CURVELETS

FELIX STÜRMER

An Evaluation of Curvlet-Based Cross-Domain Descriptors for Sketch-Based Image
Retrieval

January 2012 – version 0.1

Felix Stürmer: *Sketch-Based Image Retrieval using Curvelets*, An Evaluation of Curvlet-Based Cross-Domain Descriptors for Sketch-Based Image Retrieval, © January 2012

ABSTRACT

Short summary of the contents...

ACKNOWLEDGMENTS

acknowledgments go here...

CONTENTS

1	INTRODUCTION	1
1.1	Motivation	1
1.2	Outline	2
2	BACKGROUND & RELATED WORK	3
2.1	General Challenges of Computer Vision	3
2.1.1	The Semantic Gap	3
2.1.2	The Sensory Gap	3
2.2	Anatomy of a CBIR System	4
2.2.1	Image Acquisition	5
2.2.2	Signature Extraction	5
2.2.3	Comparison and Ranking	5
2.3	Image Transformations for Feature Extraction	5
2.3.1	The Continuous Curvelet Transform	6
2.3.2	The Fast Discrete Curvelet Transform	8
3	PROPOSED SOLUTION	13
3.1	Input Format	13
3.2	Feature Extraction	13
3.3	Distance Metric	13
4	EXPERIMENTAL RESULTS	15
5	ANALYSIS	17
6	CONCLUSION	19
	BIBLIOGRAPHY	21

LIST OF FIGURES

Figure 1	Coarse structure of a CBIR system	4
(a)	Local features	4
(b)	Global features	4
Figure 2	Signature extraction in CBIR systems	5
Figure 3	Curvelet frequency windows	7
(a)	Radial window	7
(b)	Angular window	7
(c)	Combined window	7
(d)	Complete coronisation	7
Figure 4	Discrete frequency tiling using concentric squares	9
Figure 5	Frequency tilings for USFFT and wrapping . . .	10
(a)	Sheared USFFT tiling	10
(b)	Sheared tiling for wrapping	10

LIST OF TABLES

LISTINGS

ACRONYMS

INTRODUCTION

1.1 MOTIVATION

Paragraph about increase in visual data, mobile cameras, medicine, etc...

At the core of the research into content-based image retrieval (CBIR) lies the need to be able to access the growing repositories of visual data in a convenient and efficient manner. In this context "convenient" describes the ability for the user to express the query without a complex reformulation of the intent to make it accessible to the query processor. At the same time the computational efficiency becomes more important as the amount of data to search grows. This issue becomes even more critical as the use of mobile, power-limited devices increases across many areas of application, such as autonomous vehicles or handheld augmented reality devices.

Research into text-based information retrieval has brought into existence many statistical methods to query a potentially large body of text using text as the query input. This preserves the close mapping of the intent of the user to the expression of the query and thereby makes the process accessible to users without knowledge about the internal workings of the retrieval system. Providing the means to access a large amount of visual data using a system with similar properties has turned out not to be an easy problem to solve. Using text-based querying for that purpose depends on the ability to reliably label visual data, which would require solving the general object recognition problem first [9]. To avoid that obstacle and to free the retrieval system from the requirement of translating between textual and visual information, many methods to search an image database using visual similarity have been developed.

While the goals of those systems are very similar, they differ considerably in many aspects of the processing pipeline. The query input ranges from example images over drawings to predicate describing color and shape distribution. Similarly, the structure and content of the databases and the means by which the systems query and rank the results vary significantly. This thesis focuses on evaluating a system that uses hand-drawn sketches as inputs to query databases of either full-color images or contour images. The fast discrete curvelet transform [6] is used to analyse image segments.

1.2 OUTLINE

[Chapter 2](#) presents the structure of the problem and prior solutions. The following [Chapter 3](#) proposes several variations of a particular solution using the Fast Discrete Curvelet Transform [6]. The experimental setup and its results are documented in [Chapter 4](#) and analysed in [Chapter 5](#). In [Chapter 6](#) several possible conclusions are drawn and pointers towards future research are given.

BACKGROUND & RELATED WORK

2.1 GENERAL CHALLENGES OF COMPUTER VISION

2.1.1 *The Semantic Gap*

One of the core insights of computer vision in general and content based image retrieval in specific probably is that human perception is inseparably linked to interpretation by the brain. As a human individual there is no way to directly access visual information without them having been filtered and weighted by one's personal experiences and cultural context. Therefore, when people talk about visual similarity of images, it usually includes a large degree of semantic similarity unconsciously added to the perception. The difference between that mode of perception and the current algorithmic ways to analyse visual data has been eloquently coined *the semantic gap* by [9]:

The semantic gap is the lack of coincidence between the information that one can extract from the visual data and the interpretation that the same data have for a user in a given situation.

Having had that realisation can guide the decision of a researcher or designer of such systems.

2.1.2 *The Sensory Gap*

In addition to the semantic ambiguity described above, another major obstacle of computer vision impacts a CBIR system: *the sensory gap*. This term has also been coined by [9], where it's defined as follows:

The sensory gap is the gap between the object in the world and the information in a (computational) description derived from a recording of that scene.

That terse definition includes a multitude of conditions, that can affect an image, which a CBIR system operates on:

ILLUMINATION The brightness or direction of the illumination can hide or accent edges and texture properties in the scene. Similarly, the color of the illumination influences the recorded color information in the image.

RESOLUTION The imaging resolution sets a lower limit on the size of features that can be correctly recognised by any algorithm. As in all signal processing applications, aliasing of high frequency components of the image can introduce further ambiguities. [8]

OCCLUSION Depending on the viewpoint of the recording and the composition of the scene, distinguishing parts of depicted objects may be occluded by other objects or objects may be only partially inside the recorded image.

PERSPECTIVE An object's proportions can be distorted by the imaging perspective.

An ideal CBIR system would use feature extraction and comparison methods that can account and correct for such conditions.

2.2 ANATOMY OF A CBIR SYSTEM

The inner workings of most CBIR systems can best be examined by looking at the processing pipeline each query has to go through. The coarse sequence of computational steps is almost the same in all such systems (Figure 1):

1. Acquire the image.
2. Extract the signature using a feature extraction algorithm.
3. Compare the signature to a database containing the signatures of the images to search within.
4. Rank the images by similarity using the comparison results.

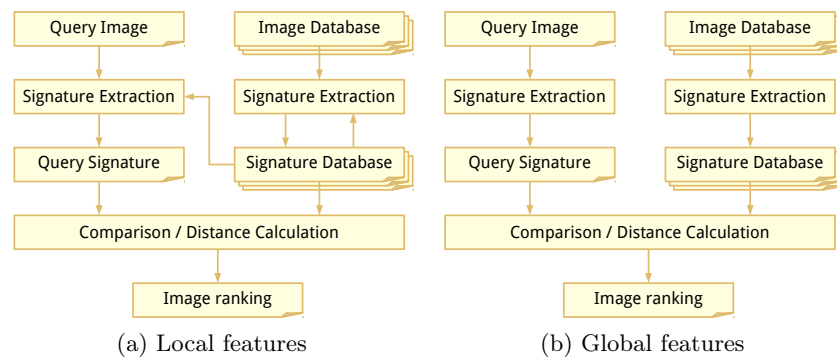


Figure 1: The processing pipeline for CBIR using both local and global features is very similar. The main difference is in the signature extraction step, in which local features are selected, weighted and/or compressed depending on the results of the signature extraction of the other images in the database.

2.2.1 Image Acquisition

The format in which the images are available to the system determines the maximum amount of information available to subsequent analysis steps.

A significant part of the preprocessing usually done after acquisition depends on the broadness of the image domain. The concept of the domain encompasses and describes the variability of many possible image parameters like illumination or composition and is therefore closely related to the sensory gap described above. The narrower the image domain is, the more assumptions the system can make about image from that domain. By their very nature, the domain of sketch based image retrieval systems is usually very broad. It contains the sketches create by the user to query the database as well as the images in the database itself, which can be of a completely different nature, e.g. photos or paintings.

Another factor usually is the accepted input format of the feature extraction algorithm. Many algorithms like SIFT [7] or SURF [2] are defined for single-channel data, but some have been specifically developed to operate on multi-channel images, like cSIFT [1].

2.2.2 Signature Extraction

TBD

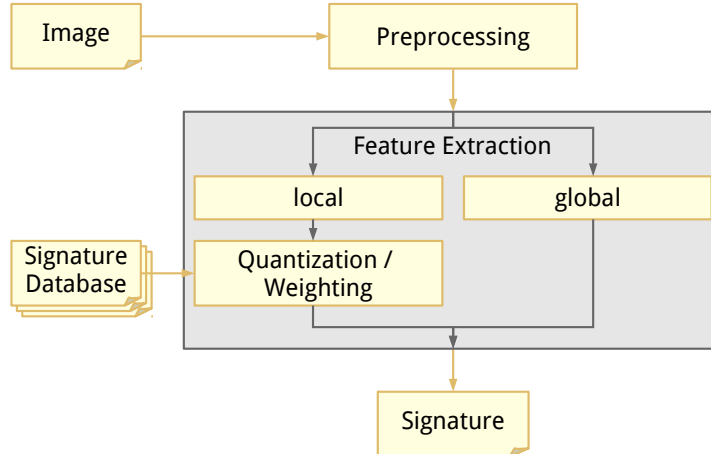


Figure 2: Signature extraction in CBIR systems

2.2.3 Comparison and Ranking

TBD

2.3 IMAGE TRANSFORMATIONS FOR FEATURE EXTRACTION

Discuss FFT, Gabor Filters, HOG, SIFT/GIST

2.3.1 *The Continuous Curvelet Transform*

The formulation of the continuous curvelet transform (CCT) by Candes and Donoho in [4] was based on Candes' previous definition and expansion of the ridgelet transform [3]. In that publication they looked at the state of research into efficient representations of edge discontinuities. They based their research on two realisations:

1. A nonadaptive approach of signal approximation can compete with many of the adaptive schemes prevalent in previous research. At the same time the non-adaptivity comes with a greatly reduced computational overhead and reduced requirements for a priori knowledge. Obtaining that knowledge in the presence of blurred or noisy data can sometimes be unfeasable.
2. Wavelet transforms can represent point singularities in a signal of up to two dimensions in a near-ideal manner, but fail to perform equally well on edges: Given a two-dimensional object in signal f , that is smooth except for discontinuities along a curve, a wavelet approximation \tilde{f}_m^W from the m largest coefficients exhibits an error of

$$\|f - \tilde{f}_m^W\|^2 \propto m^{-1}, \text{ for } m \rightarrow \infty$$

since up to $O(2^j)$ localized wavelets are needed to represent the signal along the edge. That falls short of what an approximation \tilde{f}_m^T using a series of m adapted triangles could achieve:

$$\|f - \tilde{f}_m^T\|^2 \propto m^{-2}, \text{ for } m \rightarrow \infty$$

They showed that a similarly precise approximation can be achieved by combining Candes' ridgelet analysis [3] with smart windowing functions and bandpass filters. The steps of the transformation were as follows:

1. Decomposition of the signal into subbands of scale-dependent size
2. Partitioning of each subband into squares
3. Normalisation of each square to unit scale
4. Analysis of each square in an orthonormal ridgelet system

The result was a formulation of a decomposition that matched the parabolic scaling law $\text{width} \propto \text{length}^2$ often observed in curves.

The above formulation became known as the curvelet gg transform when Candes and Donoho revised it soon after in [5]. The new version is

not dependent on ridgelets and aims to remove some shortcomings of the curvelet 99 transform, namely a simpler mathematical analysis, fewer parameters and improved efficiency regarding digital implementations, which will be described later.

The curvelet transform in \mathbb{R}^2 works by localising the curvelet waveforms in the time domain. The "mother" curvelet waveform $\varphi_j(x)$ is defined using two frequency domain windows $W(r)$, the "radial window" (Figure 3a), and $V(t)$, the "angular window" (Figure 3b). A combined frequency window U_j (Figure 3c) can then be defined as

$$U_j(r, \theta) = 2^{\frac{-3j}{4}} W(2^{-j}r) V\left(\frac{2^{\lfloor \frac{j}{2} \rfloor} \theta}{2\pi}\right).$$

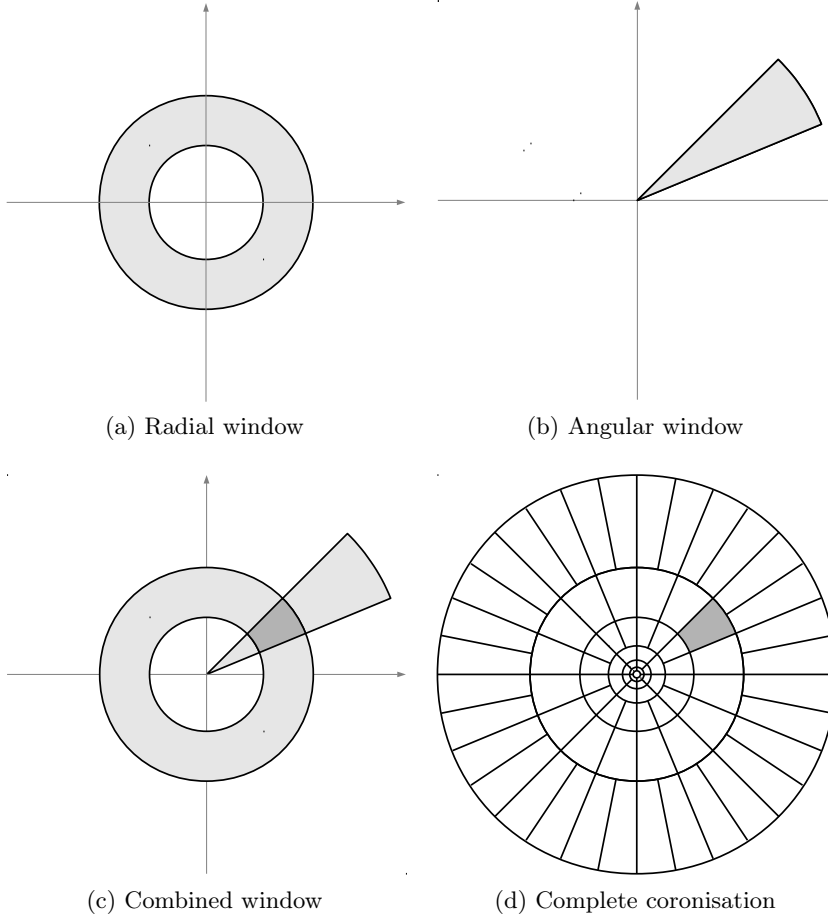


Figure 3: The window $W(2^{-j}r)$ at scale 2^j (a) is combined with the window $V(t)$ (b) to form a support wedge for the curvelet (c). The wedge roughly obeys a $\text{width} \propto \text{length}^2$ relation. (d) shows the wedge within a schema of the complete tiling in frequency domain.

The waveform φ_j can then be expressed as being the inverse Fourier transform of $\hat{\varphi}_j = U_j$ and all curvelets of a scale 2^{-j} can be derived by ROTATING φ_j by a sequence of equispaced rotation angles $\theta_l = 2\pi \cdot 2^{-\lfloor \frac{j}{2} \rfloor} \cdot l$ with $l = 0, 1, \dots$ such that $0 < \theta_l < 2\pi$ and

TRANSLATING φ_j by a sequence of offsets $\mathbf{k} = (k_1, k_2) \in \mathbb{Z}^2$:

$$\varphi_{j,k,l}(\mathbf{x}) = \varphi_j(\mathbf{R}_{\theta_l}(\mathbf{x} - \mathbf{x}_k^{(j,l)})), \quad (1)$$

where $\mathbf{x} = (x_1, x_2)$, \mathbf{R}_θ the rotation matrix for angle θ and $\mathbf{x}_k^{(j,l)} = \mathbf{R}_{\theta_l}^{-1}(k_1 \cdot 2^{-j}, k_2 \cdot 2^{-\frac{j}{2}})$.

Each curvelet coefficient $c(j, l, k)$ can then be calculated as the inner product of $f \in L^2(\mathbb{R}^2)$ and curvelet $\varphi_{j,l,k}$:

$$c(j, l, k) := \langle f, \varphi_{j,l,k} \rangle = \int_{\mathbb{R}^2} f(\mathbf{x}) \overline{\varphi_{j,l,k}(\mathbf{x})} d\mathbf{x} \quad (2)$$

As visible in figure 3d curvelets also have non-directional components at the coarsest scale, similar to those found in the wavelet transform. Those curvelets will be defined using a special low-pass filter window W_0 , which is characterized as being the remainder of the tiling not covered by the previously described radial windows:

$$|W_0(\mathbf{r})|^2 + \sum_{j \geq 0} |W(2^{-j}\mathbf{r})|^2 = 1$$

Using the window defining the coarse scale curvelet $\varphi_{j_0,k}$ via its Fourier transform is straightforward:

$$\varphi_{j_0,k}(\mathbf{x}) = \varphi_{j_0}(\mathbf{x} - 2^{-j_0}\mathbf{k}), \quad \hat{\varphi}_{j_0}(\omega) = 2^{-j_0} W_0(2^{-j_0}|\omega|), \quad (3)$$

where $\mathbf{k} = (k_1, k_2) \in \mathbb{Z}^2$.

2.3.2 The Fast Discrete Curvelet Transform

Based on the above definition of the continuous curvelet transform, a team around the authors of the original curvelet publication presented two digital, discrete implementations of the transform: the fast discrete curvelet transform (FDCT) [6]. The implementations have been described in 2D and 3D, but since this paper deals exclusively with 2D images, the explanation below will also be restricted to two dimensions.

The digital versions of the transforms operate on arrays $f[t_1, t_2]$ with $0 \leq t_1, t_2 < n$ to produce coefficients $c^D(j, l, k)$ in a way consistent with the continuous version (Equation 2):

$$c^D(j, l, k) := \sum_{0 \leq t_1, t_2 < n} f[t_1, t_2] \overline{\varphi_{j,l,k}^D[t_1, t_2]}. \quad (4)$$

Since the windows used in the continuous form are based on rotations and dyadic coroneae, they are not well suited for use with cartesian arrays. The discrete formulation substitutes them with appropriate concepts. Instead of concentric annuli, the window function W_j^D generates concentric squares "rings" using the square windows $\Phi_j(\omega_1, \omega_2) = \phi(2^{-j}\omega_1)\phi(2^{-j}\omega_2)$, with ϕ being a low-pass 1D window:

$$W_j^D(\omega) = \sqrt{\Phi_{j+1}^2(\omega) - \Phi_j^2(\omega)}, \quad j \geq 0.$$

The rotation matrix R_θ is replaced by the shear matrix S_θ to create the combined window function

$$U_{j,l}^D := W_j^D(\omega) V_j(S_{\theta_l} \omega).$$

The sequence θ_l is defined as a sequence of equispaced slopes $\tan(\theta_l) := l \cdot 2^{-\lfloor \frac{j}{2} \rfloor}$ with $l = -2^{\lfloor \frac{j}{2} \rfloor}, \dots, 2^{\lfloor \frac{j}{2} \rfloor} - 1$. The construction of the corner windows are Figure 4 shows a tiling of all $U_{j,l}^D$.

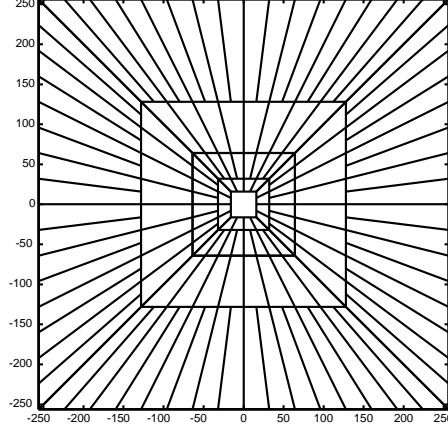


Figure 4: Discrete frequency tiling using concentric squares

2.3.2.1 FDCT using unequipped FFTs

The first implementation of the discrete curvelet transform transfers the input array $f[t_1, t_2]$, $0 \leq t_1, t_2 < n$ into the Fourier domain to obtain $\hat{f}[n_1, n_2]$:

$$\hat{f}[n_1, n_2] = \sum_{t_1, t_2=0}^{n-1} f[t_1, t_2] e^{-\frac{i2\pi(n_1 t_1 + n_2 t_2)}{n}}, \quad -\frac{n}{2} \leq n_1, n_2 < \frac{n}{2}$$

The obtained Fourier samples need to be interpolated for each pair of scale j and angle l to match the grid of the sheared support window $U_j^D[n_1, n_2]$. The authors achieve this by resampling \hat{f} on the grid implied by the sheared window for each angle via a series of 1D fast Fourier transforms. These transforms represent a polynomial interpolation of each "column" of the parallelogram P_j containing the sheared window (Figure 5a), that can be computed with a $O(n^2 \log n)$ complexity in a sufficiently exact approximation.

This yields an object $\hat{f}[n_1, n_2 - n_1 \tan \theta_l]$ for $(n_1, n_2) \in P_j$, that can be multiplied with the window U_j^D described above in order to create a localized "wedge" with the orientation θ_l :

$$f_{j,l}^D[n_1, n_2] = \hat{f}[n_1, n_2 - n_1 \tan \theta_l] U_j^D[n_1, n_2]$$

The discrete curvelet coefficients $c^D(j, l, k)$ can then be calculated by applying the inverse 2d Fourier transform:

$$c^D(j, k, l) = \sum_{n_1, n_2 \in P_j} \hat{f}[n_1, n_2 - n_1 \tan \theta_l] U_j^D[n_1, n_2] e^{i2\pi(\frac{k_1 n_1}{L_{1,j}} + \frac{k_2 n_2}{L_{2,j}})},$$

in which $L_{1,j}$ and $L_{2,j}$ are the length and width of the rectangle supporting U_j^D .

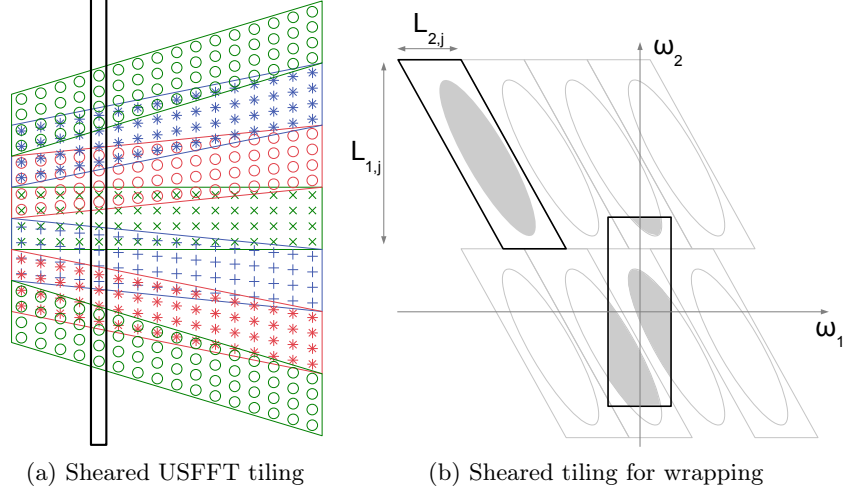


Figure 5: (a) illustrates the respective grid for each parallelogram containing the sheared support windows of the "east quadrant". The box highlights one of the columns, that represent one of the 1D polynomial interpolation problems solved for resampling. In (b) the parallelogram $P_{j,l}$ is shown on top of the tilted tiling of a curvelet in frequency domain. Due to the periodicity the rectangle in the center contains the same curvelet, but has a much smaller axis aligned bounding box for the FFT to operate on.

2.3.2.2 FDCT using wrapping

As before, FDCT using wrapping first calculates $\hat{f}[n_1, n_2]$ with $-\frac{n}{2} \leq n_1, n_2 < \frac{n}{2}$ as the Fourier transform of the input f . The sample is then localized by multiplying it with a window $U_{j,l}^D$ for each angle j and scale l :

$$d_{j,l}[n_1, n_2] = U_{j,l}^D \hat{f}[n_1, n_2]$$

To avoid the computationally costly interpolation step required in the USFFT approach, this method keeps the rectangular grid of the input signal. Because an axis-aligned bounding box of the window U_j^D in Fourier domain cannot maintain the $\text{width} \propto \text{length}^2$ proportions of the window, applying an inverse Fourier transform on such a bounding box in general would lead to significant oversampling of the coefficients and thereby increase the memory requirements for fine scale curvelets

beyond that of the the USFFT approach. In order to circumvent that, the authors utilize the periodic nature of the Fourier transform and propose generating a periodically wrapped version of the fourier samples. For $P_{j,l}$ as the bounding parallelogram of $U_{j,l}^D$, $L_{1,j}$ and $L_{2,j}$ are the period lengths by which to translate $P_{j,l}$ in the horizontal and vertical direction to produce a suitable tiling for each orientation θ_l (Figure 5). Thus, the wrapped, localized data are

$$f_{j,l}^D[n_1, n_2] = Wd_{j,l}[n_1, n_2] = \sum_{m_1 \in \mathbb{Z}} \sum_{m_2 \in \mathbb{Z}} d_{j,l}[n_1 + m_1 L_{1,j}, n_2 + m_2 L_{2,j}]$$

with $0 \leq n_1 < L_{1,j}$ and $0 \leq n_2 < L_{2,j}$, which gives a rectangle of size $L_{1,j}$ times $L_{2,j}$.

Again, the discrete curvelet coefficients $c^D(j, l, k)$ can then be collected using the inverse 2d Fourier transform:

$$c^D(j, k, l) = \sum_{n_1, n_2 \in P_j} f_{j,l}^D[n_1, n_2] e^{i2\pi(\frac{k_1 n_1}{L_{1,j}} + \frac{k_2 n_2}{L_{2,j}})}.$$

PROPOSED SOLUTION

Proposed solution goes here. . .

3.1 INPUT FORMAT

- Luma component (Y') of $Y'UV$ representation
- Gradient magnitude of Sobel operator of luma component
- Canny edge map of luma component
- gPb

3.2 FEATURE EXTRACTION

- Global features: mean and standard deviation
- Local features: visual words via k-means clustering
- great comparison of sampling for k-means clustered vws [nowako6]

3.3 DISTANCE METRIC

- Euclidean Distance
- cosine distance?
- EMD?

EXPERIMENTAL RESULTS

Experimental results go here. . .

ANALYSIS

Analysis goes here. . .

CONCLUSION

Conclusion goes here. . .

BIBLIOGRAPHY

- [1] A.E. Abdel-Hakim and A.A. Farag. “CSIFT: A SIFT Descriptor with Color Invariant Characteristics.” In: *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. Vol. 2. 2006, pp. 1978–1983. DOI: [10.1109/CVPR.2006.95](https://doi.org/10.1109/CVPR.2006.95).
- [2] H. Bay et al. “Speeded-up robust features (SURF).” In: *Computer Vision and Image Understanding* 110.3 (2008), 346–359. URL: <http://www.sciencedirect.com/science/article/pii/S1077314207001555>.
- [3] E. J. Candes. “Ridgelets: theory and applications.” PhD thesis. Stanford University, 1998. URL: <http://www-stat.stanford.edu/~candes/papers/thesis.ps>.
- [4] E. J. Candes and D. L. Donoho. *Curvelets: A surprisingly effective nonadaptive representation for objects with edges*. Tech. rep. DTIC Document, 2000. URL: <http://oai.dtic.mil/oai/oai?verb=getRecord&metadataPrefix=html&identifier=ADP011978>.
- [5] E. J. Candes and D. L. Donoho. “New tight frames of curvelets and optimal representations of objects with piecewise C² singularities.” In: *Communications on pure and applied mathematics* 57.2 (2004), 219–266. URL: <http://onlinelibrary.wiley.com/doi/10.1002/cpa.10116/abstract>.
- [6] E. Candes et al. “Fast discrete curvelet transforms.” In: *Multiscale modeling and simulation* 5.3 (2006), 861–899.
- [7] D. G. Lowe. “Object recognition from local scale-invariant features.” In: *Computer Vision, 1999. The Proceedings of the Seventh IEEE International Conference on*. Vol. 2. 1999, 1150–1157. URL: http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=790410.
- [8] C. E. Shannon. “Communication in the presence of noise.” In: *Proceedings of the IEEE* 86.2 (1998), 447–457. URL: http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=659497.
- [9] A. W.M Smeulders et al. “Content-based image retrieval at the end of the early years.” In: *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 22.12 (2000), 1349–1380.

COLOPHON

This document was typeset using the typographical look-and-feel `classicthesis` developed by André Miede. The style was inspired by Robert Bringhurst’s seminal book on typography “*The Elements of Typographic Style*”. `classicthesis` is available for both \LaTeX and \L\X :

<http://code.google.com/p/classicthesis/>

Happy users of `classicthesis` usually send a real postcard to the author, a collection of postcards received so far is featured here:

<http://postcards.miede.de/>

DECLARATION

Put your declaration here.

Berlin, January 2012

Felix Stürmer