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

# ABSTRACT

Short summary of the contents. . .

# ACKNOWLEDGMENTS

acknowledgments go here...

# CONTENTS

[ August 9, 2012 at 15:25 – `classicthesis` version 0.1 ]

## LIST OF FIGURES

## LIST OF TABLES

## LISTINGS

## ACRONYMS

[ August 9, 2012 at 15:25 – classicthesis version 0.1 ]

# 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 [6]. 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 [4] 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 [4]. 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 unconciously 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 [6]:

> 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 [6], 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. [5]

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.

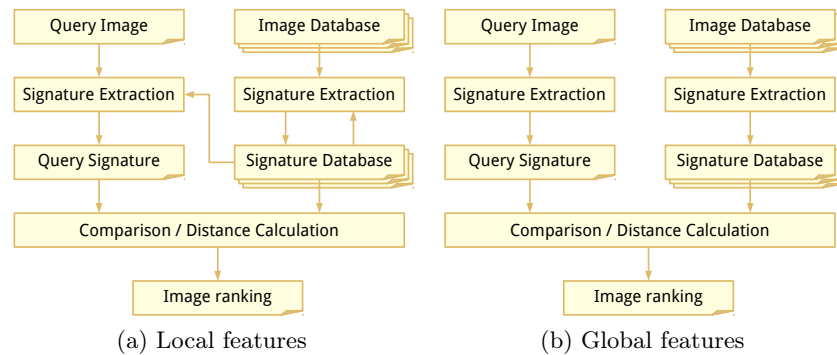4. Rank the images by similarity using the comparison results.



(a) Local features                (b) Global features

Figure 1: Coarse structure of a CBIR system

### 2.2.1 Image Aquisition
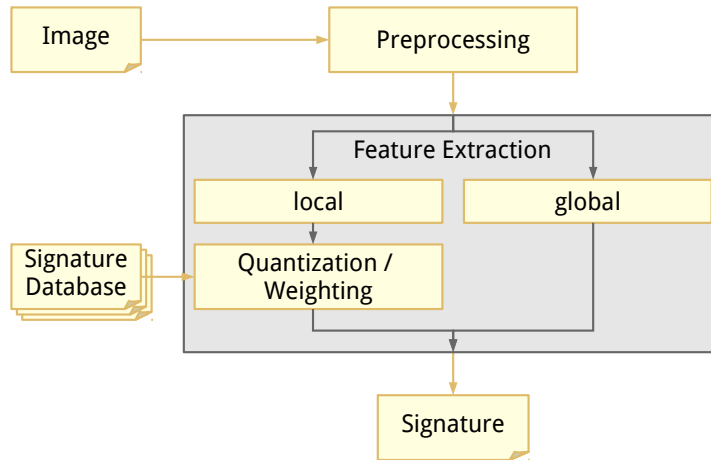
TBD

### 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 [2] was based on Candes' previous definition and expansion of the ridgelet transform [1]. 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

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

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:

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

They showed that a similarly precise approximation can be achieved by combining Candes' ridgelet analysis [1] 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 $\texttt{width} \propto \texttt{length}^2$ often observed in curves.

The above formulation became known as the curvelet 99 transform when Candes and Donoho revised it soon after in [3]. 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

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

$\varphi_j$ 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 $\varphi_j$ by a sequence of equispaced rotation angles $\theta_l = 2\pi \cdot 2^{-\lfloor \frac{j}{2} \rfloor} \cdot l$ with $l = 0, 1, \ldots$ such that $0 < \theta_l < 2\pi$ and

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

$$\varphi_{j,k,l}(x) = \varphi_j(R_{\theta_l}(x - x_k^{(j,l)})), \tag{1}$$

where $x = (x_1, x_2)$, $R_\theta$ the rotation matrix for angle $\theta$ and $x_k^{(j,l)} = R_{\theta_l}^{-1}(k_1 \cdot 2^{-j}, k_2 \cdot 2^{-\frac{j}{2}})$.

(a) Radial window

(b) Angular window
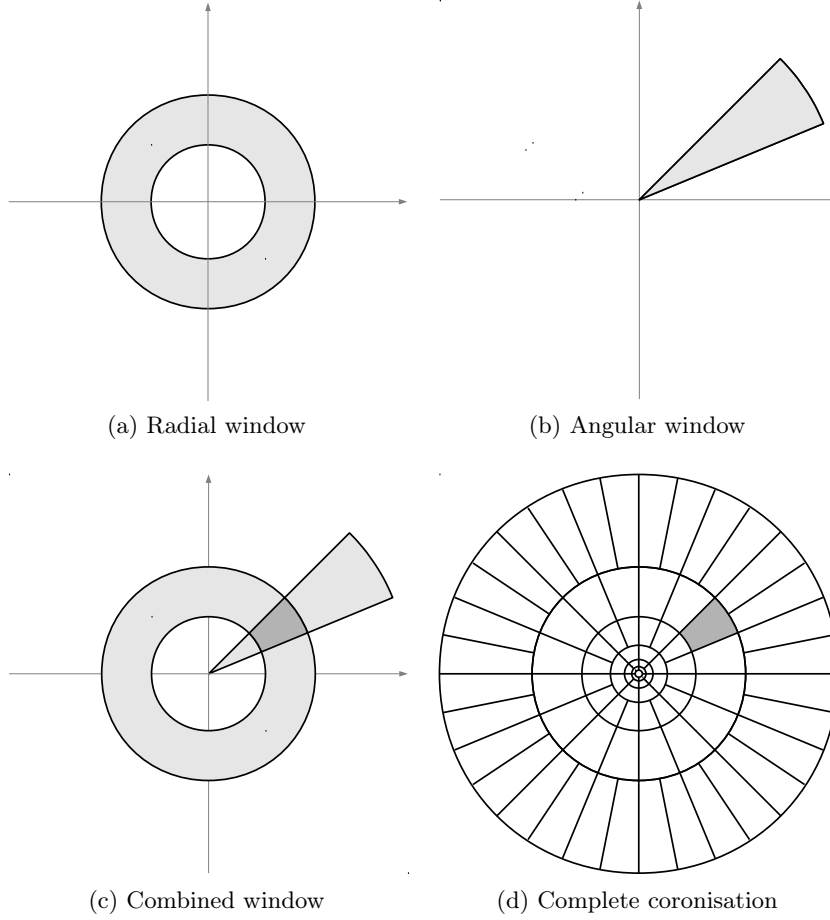
(c) Combined window

(d) Complete coronisation

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 width $\propto$ length$^2$ relation. (d) shows the wedge within a schema of the complete tiling in frequency domain.

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 $\varphi_{j,l,k}$:

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

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:

$$|W_0(r)^2| + \sum_{j \geqslant 0} |W_({2^{-j}r)}|^2 = 1$$

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

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

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

### 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 [4]. 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 $0 \leqslant t_1, t_2 < n$ to produce coefficients $c^D(j, l, k)$ in a way consistent with the continuous version (Equation 2):

$$c^D(j, l, k) := \sum_{0 \leqslant t_1, t_2 < n} f[t_1, t_2] \overline{\varphi^D_{j,l,k}[t_1, t_2]}. \tag{4}$$

# 3

## 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 [nowak06]

### 3.3 DISTANCE METRIC

- Euclidean Distance
- cosine distance?
- EMD?

9

# EXPERIMENTAL RESULTS

4

Experimental results go here...

# 5

## ANALYSIS

Analysis goes here. . .

# CONCLUSION

6

Conclusion goes here...

## BIBLIOGRAPHY

[1] E. J. Candes. "Ridgelets: theory and applications." PhD thesis. Stanford University, 1998. URL: http://www-stat.stanford.edu/~candes/papers/thesis.ps.

[2] 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.

[3] 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.

[4] E. Candes et al. "Fast discrete curvelet transforms." In: *Multiscale modeling and simulation* 5.3 (2006), 861–899.

[5] 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.

[6] 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.

[ August 9, 2012 at 15:25 – `classicthesis` version 0.1 ]

# DECLARATION

Put your declaration here.

*Berlin, January 2012*

_____

Felix Stürmer