

# Argument for Scene Categorization with Image Amplitude Spectra

## Purpose

Although many researches on scene recognition have shown that semantic information is mainly conveyed by the phase spectrum [1, 2], most of psychophysical works [3-5] paradoxically use the amplitude spectrum as the principal component of a natural scene. The aim of our work was to test whether amplitude without phase spectrum data contain sufficient information to classify natural scenes. Ninety-two noisy natural scenes belonging to four categories (Beach, City, Forest and Indoor Scene) were presented both to human subjects and to a classification model which reproduces early human visual processes (retina, V1 cortical cells). The shapes of the noise's amplitude spectrum used were those of cities ( $N_c$ ), of beaches ( $N_b$ ), of generic natural scenes [6], or white noise. Behavioral and simulation data show that scene categorization was easier for pictures composed of a congruent noise than for incongruent ones. This leads us to think that amplitude without phase spectrum could be sufficient to make natural scene classification.

## Amplitude spectrum

The categorical noise contains only information relative to the scene category in the amplitude spectrum. Then, all (other) things being equal, if performance on scene categorization varies as a function of the "category" of the noise, then the amplitude spectrum carries some useful category information processed by the visual system.

## Stimuli

92 stimuli were built from 28 natural scene grey-level pictures and noises from four categories. 7 pictures ( $I_b$  from four categories (Beach  $I_b$ , City  $I_c$ , Forest  $I_f$  and Indoor Scene  $I_i$ ) were chosen according to their amplitude spectrum's shape: vertical for beach, horizontal for city, both for indoor and isotropic for forest [8, 9]. Four classes of noise ( $N_m$ ) were built: horizontal  $N_m$ , vertical  $N_m$ , white  $N_m$  and natural  $N_m$ . The shapes of noises were obtained (cf. equation 1) by multiplying white noise  $w(x)$  with the mean amplitude spectra  $mm(x)$  of city and beach (for  $N_c$  and  $N_b$ ) or with 1/f spectra (for  $N_m$ ) [7].

$$N_m(x) = \frac{DFT^{-1} [ DFT [ w(x) ] * DFT [ m_m(x) ] ]}{\varepsilon(N_m)} * \varepsilon(I_m) \quad (1)$$

$$\varepsilon [ f(x) ] = \sum_x | DFT [ f(x) ] |^2 \quad (2)$$

Each noise is normalized (cf. equation 2), and a coefficient  $k$  is applied in order to ensure a sufficient masking. On the basis of a pre-test experiment, we choose 1.75 as coefficient value, which resulted in a categorization performance superior to 80% accuracy.

When adding noise to natural scenes (cf. equation 3), their respective amplitude spectra could be neutral, congruent (Beach with vertical noise or City with horizontal noise) or incongruent according to the scene category (cf. figure 1).

$$I_{pm}(x) = I_p(x) + k * N_m(x) \quad (3)$$

As  $w(x)$  is a white noise, its module is constant and its phase random,

$$N_m(x) = |M_m(x)| e^{j \cdot rand(x)} \quad (4)$$

If ones swaps the phase spectrum of a scene, it contains no more image specific information and therefore the picture could not be recognisable [2]. We could assess that if the phase spectrum of a signal is multiplied with a random signal, the resulting signal is random:

$$N_m(x) = |M_m(x)| e^{j \cdot rand(x)} \quad (5)$$

## Loci of Information

## Gabor Model

Former studies have shown that the global distribution of the local dominant orientations appears to be a powerful feature for discriminating between four semantic categories of real world scenes [10].

### Retinal pre-processing

The retinal photoreceptors make a space-time high-pass filtering after an adaptive compression process. This results in a contrast equalization of the image and spectral whitening which compensates for the 1/f image amplitude spectrum [11].

### Cortical filtering

In area V1 of the visual cortex, the retinal image is decomposed by the filtering of cortical neurons, which are sensitive to various spatial frequency bands and various orientations of the stimuli. Here, we aim at categorizing and not describing scenes, so we simulate complex cells which are invariant to object position in the scene. These cells, described by the means of Gabor wavelets, provide the local energy of images. Images are filtered by Gabor wavelets into 7 frequency bands and 7 different orientations. According to the physiological data about the visual cells [12], the relative radial bandwidth of the Gabor filters is fixed at 1 octave. So, each image is analyzed with a bank of 49 Gabor filters, the output energies of which provide a point in a 49-dimensions space.

## Results

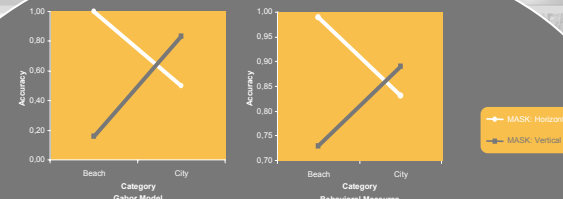


Figure 2 : Percentage of accuracy as a function of the scene Category (Beach vs City) and the amplitude spectrum characteristics of Masks (Horizontal vs Vertical). (The scales of plots differ).

## Participants

Eleven students from the National Polytechnic Institute of Grenoble (ages between 18 and 25, with normal or corrected vision) were volunteers to participate in the experiments.

## Procedure

The participants sat in front of the screen at a distance of 1.2 meters. Stimuli subtended a visual angle of 9.7°. Each picture was displayed for 20 ms; it was preceded by a black central fixation dot for 1.5 s, and followed by a white screen for 1.5 s. Participants were to press a button as quickly and as accurately as possible if the picture was a beach or city, and not to do so for the other picture categories ('Go - no Go' procedure). Reaction times were recorded with a vocal key and responses were keyed to disk by the experimenter. No feedback was given. Each experiment started with 2 practice sets. The order of trials was random.

## Behavioral Measures

## Gabor Model

In order to measure our model accuracy, we compute the Euclidean distance between the amplitude spectrum of each noisy image (i.e. Beach and City with  $N_b$  and  $N_c$  noises) and the mean amplitude spectra of beach and city categories. Images are associated to the closest category. Accuracy is the ratio of good associations to the cardinal of each category. (Figure 2a)

## Results

### Behavioral Measures

A 2 x 2 ANOVA (Mask x Category), carried out on accuracy, showed no effect of Category ( $F(1,10) < 1$ ), an effect of Mask ( $F(1,10)=6.33, p<0.03$ ): Participants responded more accurately for horizontal mask ( $N_b$ ) than for vertical mask ( $N_c$ ), and a Mask x Category interaction ( $F(1,10)=10.14, p<0.01$ ): Participants responded more accurately when the amplitude spectra of the mask and the category were congruent ( $I_b$  and  $I_c$  conditions) than incongruent ( $I_b$  and  $I_b$  conditions) ( $F(1,10)=12.025, p<0.01$ ). A 2 x 2 ANOVA (Mask x Category), carried out on correct response times, leads to similar results. (Figure 2b)

## Discussion

With this work, we have shown that image amplitude spectrum is sufficient to make some natural scene clusters. So, shaping noise is an efficient paradigm in order to investigate low-level vision processes. Moreover, these results provide us a validation of our model of scene categorization based on the global distribution of image energy.

There are still two main problems in the experiment:

- Noise shapes could be changed to render them more orthogonal to the shape of semantical category amplitude spectra [7, 10, 11]
- The signal to noise ratio (SNR) of the stimuli computed in the spatial domain can also explain our results (even if it is equal to 1 in the frequency domain). Now, the SNR cannot account for the structural organisation of a natural scene. Equalizing the histogram produces a SNR  $\neq 1$  without modifying the structure of the picture. However if one swaps each pixel randomly the SNR will be equal to 1 and the picture become noise



Sponsors

Alan Chauvin\*, Nathalie Guyader\*, Christian Marendaz\*, Jeanny Héroult\*,  
 \*Laboratoire de Psychologie et NeuroCognition - UMR 5105, UPMF BP 47, 38040 Grenoble, France  
 \*Laboratoire des Images et des Signaux - UMR 5083, INPG, 46 Av. Félix Viallet 38031 Grenoble Cedex,

1. Ogden, J.M. and Lim, J.S. The importance of phase for objects: Proc. IEEE, 1981, 69, pp. 552-551.
2. Morgan, M.J., Ross, J., and Hayes, A. The relative importance of local phase and local amplitude in patchwise image reconstruction. Biological Cybernetics, 1991, 66, p. 113-119.
3. Földiák, D. and Brady, N. Visual sensitivity, blur and the sources of variability in the amplitude spectra of natural scenes. Vision Res, 1997, 37(23), p. 3367-83.
4. Ruderman, D.L. Origins of scaling in natural images. Vision Res, 1997, 37(23), p. 3331-50.
5. Torralba, J. and Sintoniz, Y. Bandwidth of contrast in natural images: Relations to detectability of changes in the amplitude spectra. Vision Res, 1997, 37(23), p. 3203-15.
6. Van der Schaaf, A., and Van Hateren, J.H. Modeling the power spectra of natural images: statistics and information. Vision Res, 1996, 36(17), p. 2239-7.
7. Guérou-Dupuy, A. and Osta, A. Classifications of scene photographs from local orientations features. Pattern Recognition Letters, 2009, 29, p. 1939-1940.
8. Héroult, J. De la rétroanalyse aux caractéristiques neuronales, in Les Systèmes de Vision, J.M. Jonjat, Editeur, 2001, Hermès.
9. De Vries, S.M., and de Vries, K.J. Spatial Vision, 1983, Oxford: Oxford University Press.
10. Osta, A. and Torralba, A.B. Modeling the Shape of the Scene: A Holistic Representation of the Spatial Envelope. International Journal of Computer Vision, 2001, 43(3), p. 145-176.
11. Guyader, N. and Héroult, J. Représentation espace-féquence pour la catégorisation d'images. in GRETSI, 2001, Toulouse.