# A Signal Processing Approach to Modeling Low-level Vision, and Applications

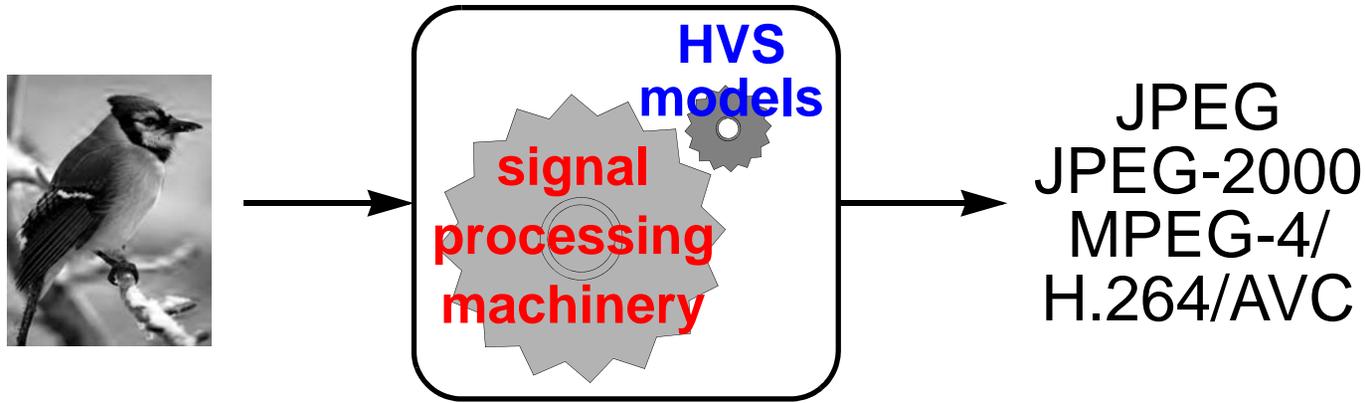## Sheila S. Hemami

Department of Electrical & Computer Engineering
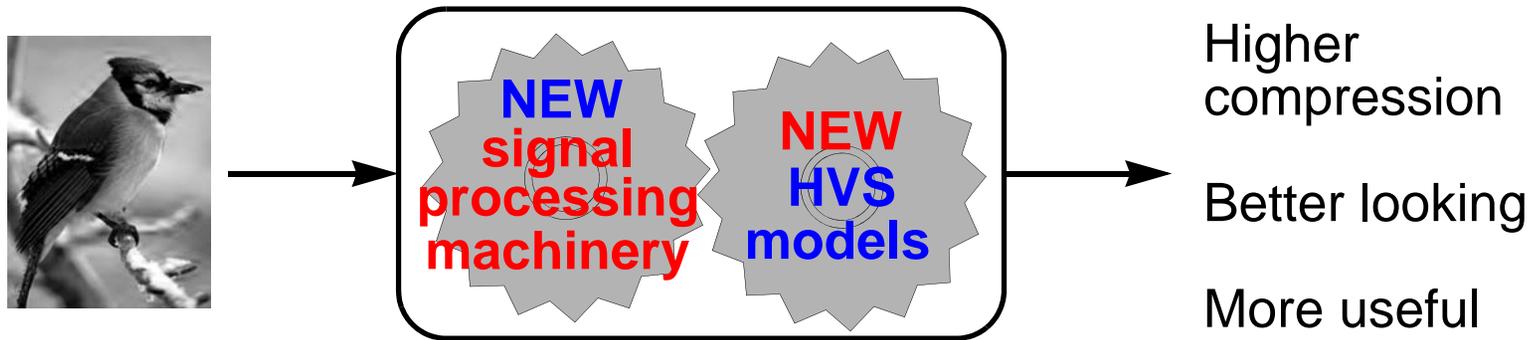Northeastern University

## Visual Communications Lab

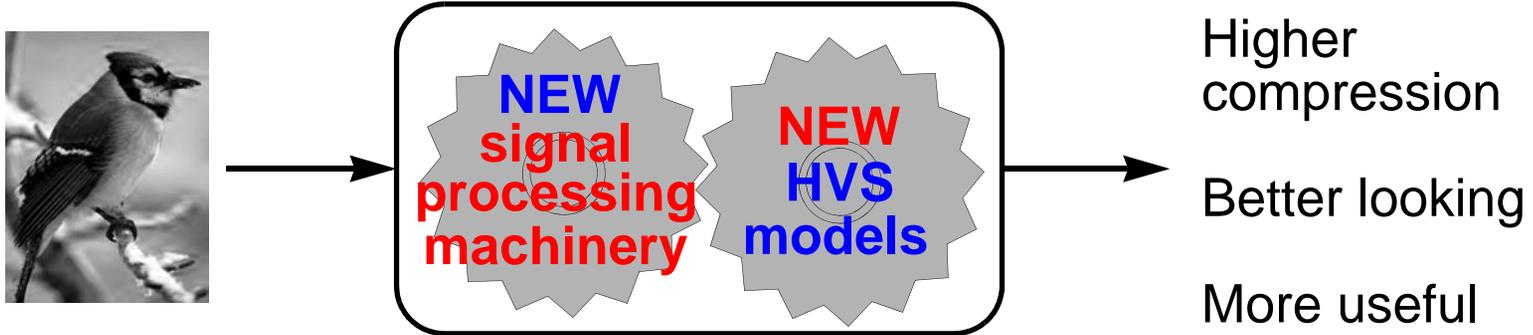- A model for current state-of-the-art techniques:



- This model works well when *transmission resources are not limited* (bandwidth, QoS, etc.).

- When resources become scarce, *every bit counts*. The SP machinery starts to break at low rates.

Higher compression

Better looking

More useful

- Develop *more appropriate HVS models* suitable for image applications via strategic psychophysical experimentation.

- Develop signal processing theory and practice to exploit this HVS characterization.

- Develop *more appropriate HVS models* suitable for image applications via strategic psychophysical experimentation.

- Develop signal processing theory and practice to exploit this HVS characterization.

- End goal: an image processing system incorporating a model which exhibits better performance than if the model is not used.
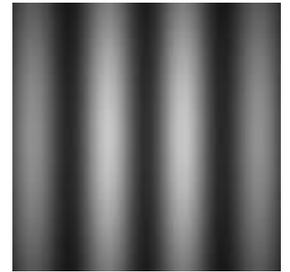
# Outline

- Three "classical" psychophysics results/HVS characterizations.

- Wavelets, the multichannel model, and images.

- Characterizing the HVS using natural images.

- Some SP strategies and applications to compression which exploit our characterization.
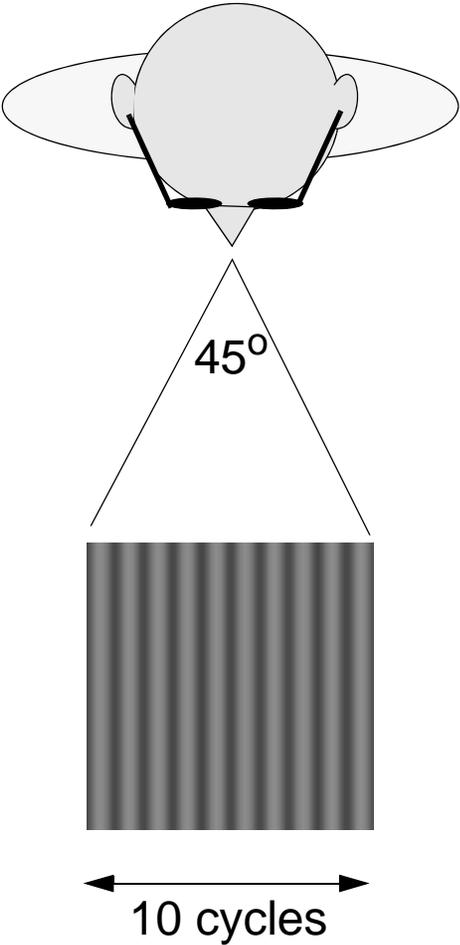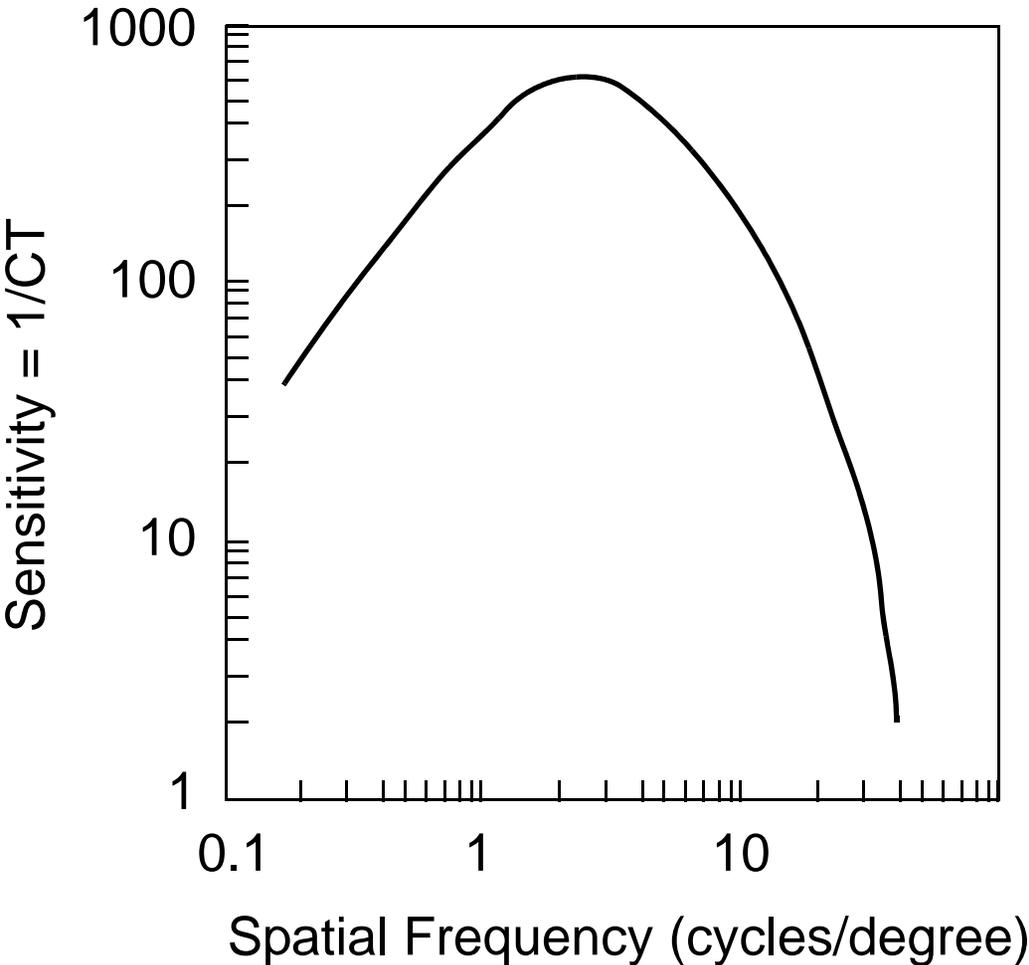
- Current work: image utility.

Experiments with sinusoidal gratings
yield the following:

1. The human *contrast sensitivity function* (CSF)

1. The CSF is *not orientation specific*.

2. The CSF is *for simple gratings only*.

3. The CSF represents *subthreshold* perception.

The HVS consists of *channels*, each tuned to range of spatial frequencies and orientations.

1. The CSF is *not orientation specific*.

2. The CSF is *for simple gratings only*.

3. The CSF represents *subthreshold* perception.

Sensitivity = 1/CT

Spatial Frequency (cycles/degree)

1000
100
10
1

0.1
1
10

Reference grating

Spatial Frequency (c/deg)

Two gratings at different frequencies have equal perceived contrast at equal physical contrast as they become increasingly suprathreshold.

Experiments with sinusoidal gratings
yield the following:



1. The human *contrast sensitivity function* (CSF) —
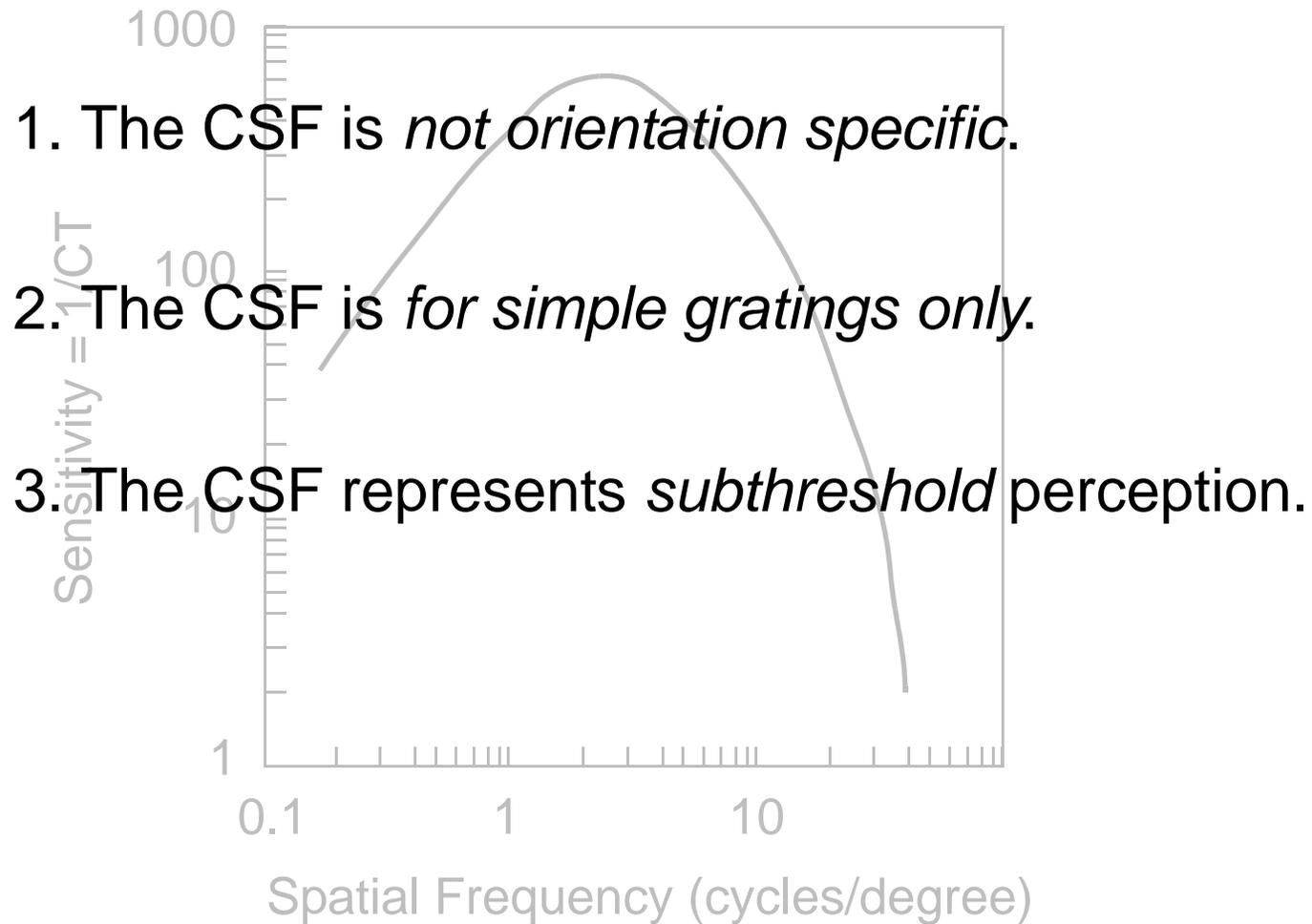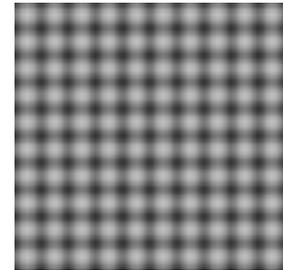   the HVS has a low-pass response at the
   detection threshold, becoming flat as gratings
   become more  visible.

2. *Summation*

If this stimulus has contrast threshold

$$CT_A$$

...and this stimulus has contrast threshold

$$CT_B$$

Then what is the contrast threshold of this stimulus?

$$CT_{A+B} = \; ?$$

If this stimulus has contrast threshold

$CT_A$

...and this stimulus has contrast threshold

$CT_B$

Then what is the contrast threshold of this stimulus?

$CT_{A+B} = ?$

- For the compound stimuli to be as detectable as either of the individual components,

$$(C_A / CT_A)^\beta + (C_B / CT_B)^\beta = 1$$

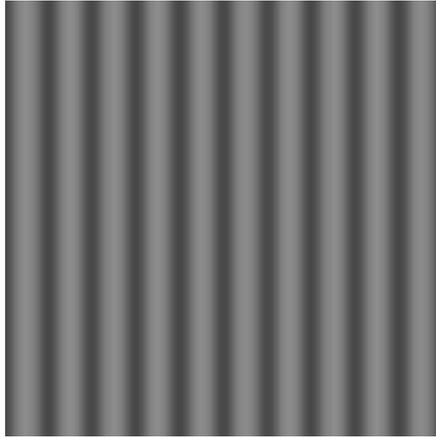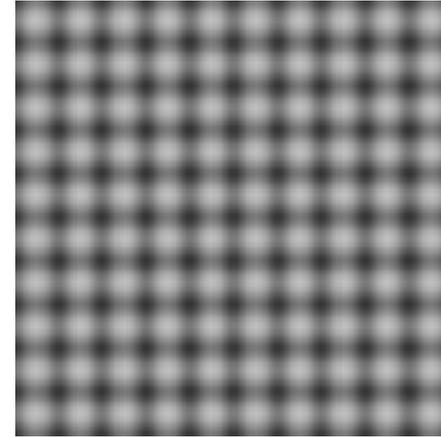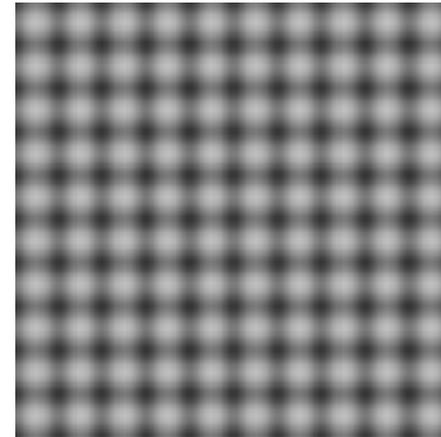- For sinusoidal components, $\beta \in [2, 4]$.

Experiments with sinusoidal gratings
yield the following:

1. The human *contrast sensitivity function* (CSF) — the HVS has a low-pass response at the detection threshold, becoming flat as gratings become more visible.

2. *Summation* — The contrast threshold for a given sinusoid is 40% lower when it is shown simultaneously with another, different sinusoid.

Experiments with sinusoidal gratings
yield the following:

1. The human *contrast sensitivity function* (CSF) —
the HVS has a low-pass response at the
detection threshold, becoming flat as gratings
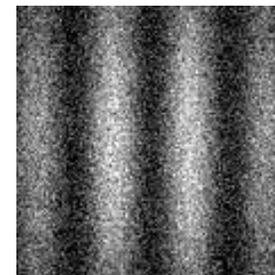become more visible.

2. *Summation* — The contrast threshold for a given sinusoid is
40% lower when it is shown simultaneously with another,
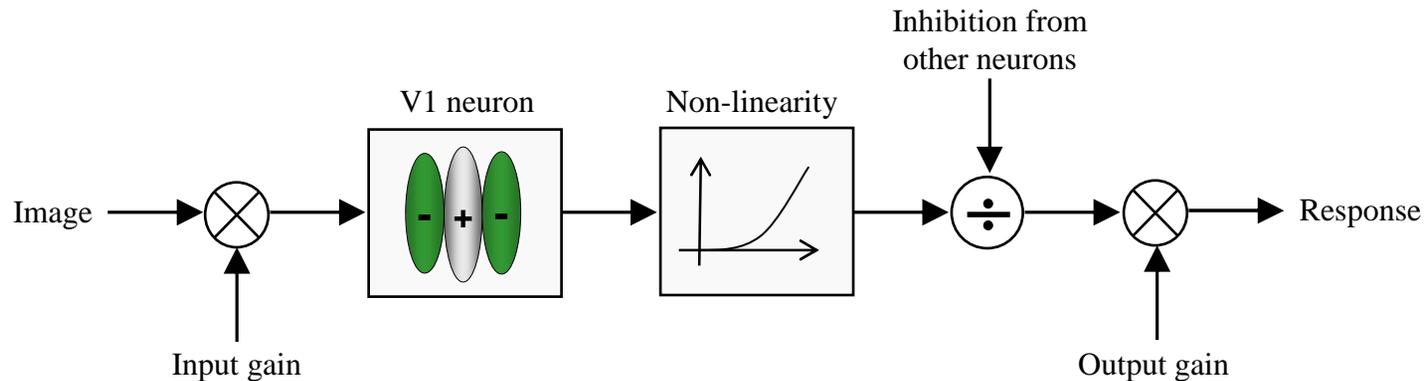different sinusoid.

3. The standard gain control model for *masking* describes how
thresholds are impacted based on surrounding image content.

# Standard Gain Control Model (Masking)



Neural response:
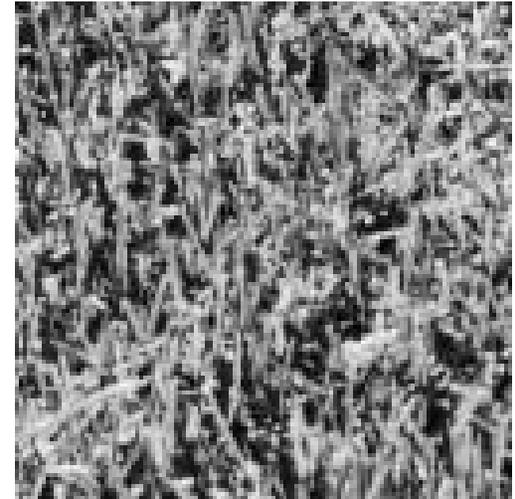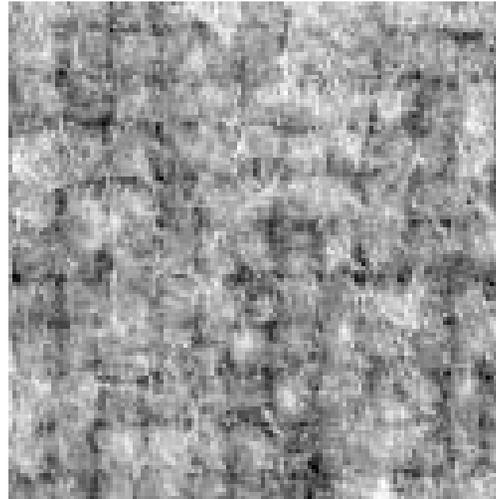
$$r(x, f, \theta) = \frac{w(x, f, \theta)^p}{b^q + \sum_{(x, f, \theta) \in S} w(x, f, \theta)^q}$$

$w$ = e.g., wavelet coefficient at location $x$, frequency $f$, orientation $\theta$

Usually $p \approx 2$, $q \approx 2$ — effectively variance!

# Standard Gain Control Applied to *Textures*



The standard visual masking model predicts the masking
elevations well *for homogeneous textures*.

# The Signal Processor's Question

Should the 3 classical psychophysical results, based on sinusoidal gratings be directly applied to processing images?

- Images are the superposition of many sinusoidal components.

- Images provide a very sophisticated "mask" to any distortions introduced by compression.

- Arbitrary image patches are not necessarily homogeneous textures.

- [Images have higher-level meaning to observers.]

Should the 3 classical psychophysical results, based on sinusoidal gratings be directly applied to processing images?
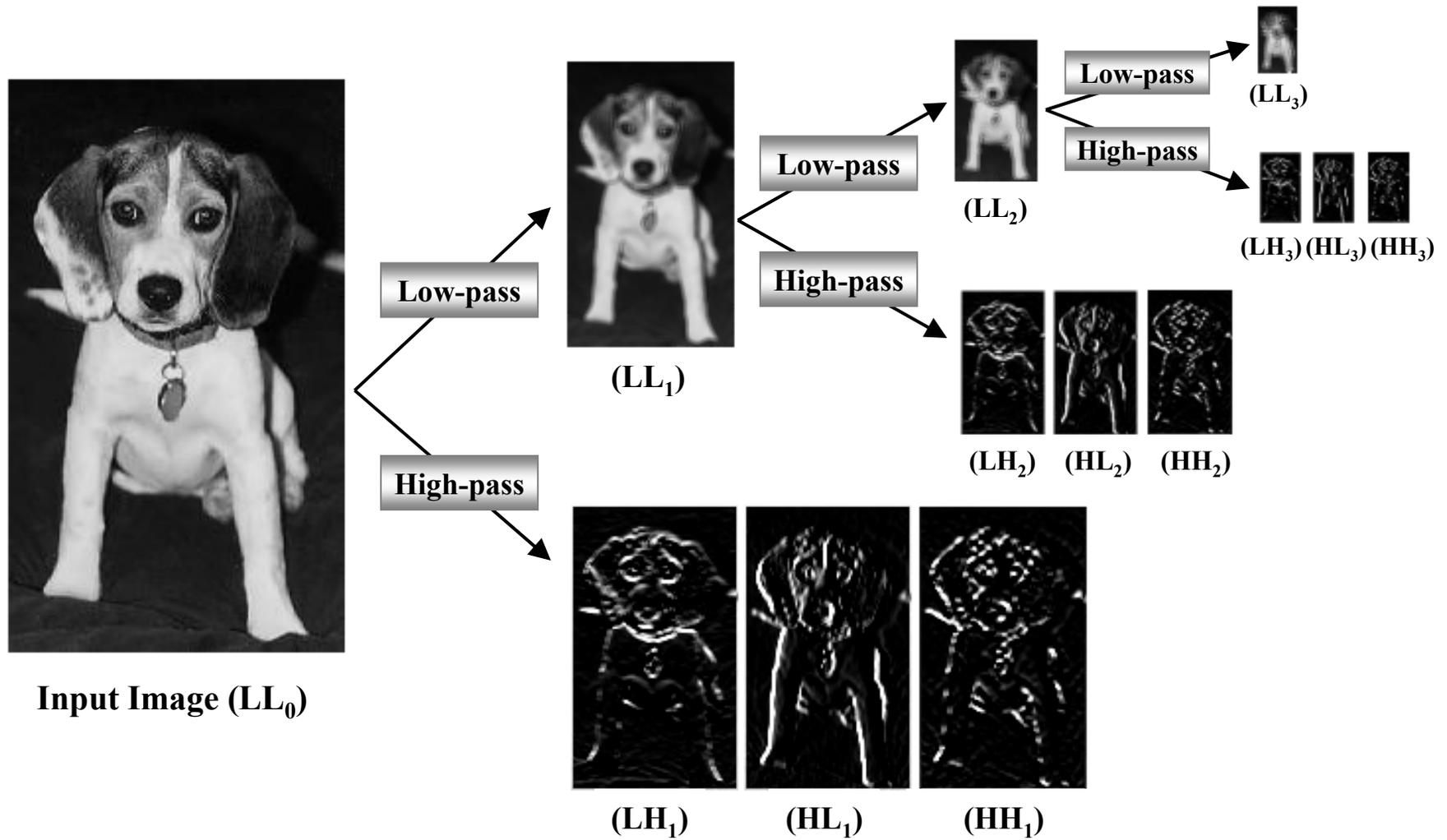
# NO

# Questions

Using realistic maskers (images) and realistic stimuli (bandlimited, correlated quantization noise)...

- What are the visibility thresholds for *quantization distortions as occur in natural images*? (CSF without and with masking)

- How are distortions from multiple quantized subbands perceived? (Summation)

- Can we predict visibility thresholds from *local natural image characteristics*? (Masking)

- [How should higher-level processing (i.e., the task) impact any necessary signal processing?]
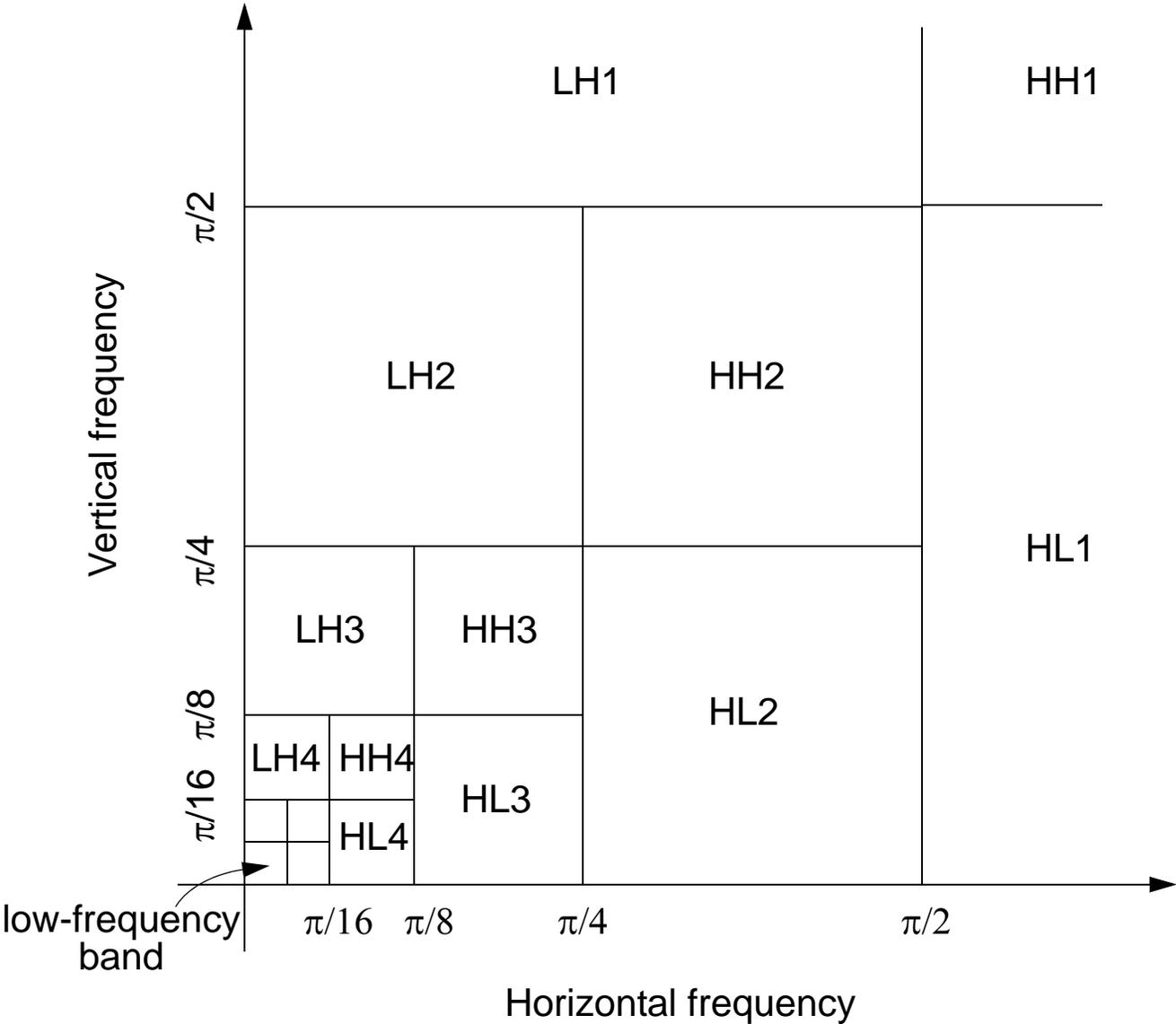
# Outline

- Three "classical" psychophysics results/HVS characterizations.

- <span style="color:red">Our image coding framework: wavelets, the multichannel model, and digital images.</span>

- Characterizing the HVS using natural images.

- Some SP strategies and applications to compression which exploit our characterization.
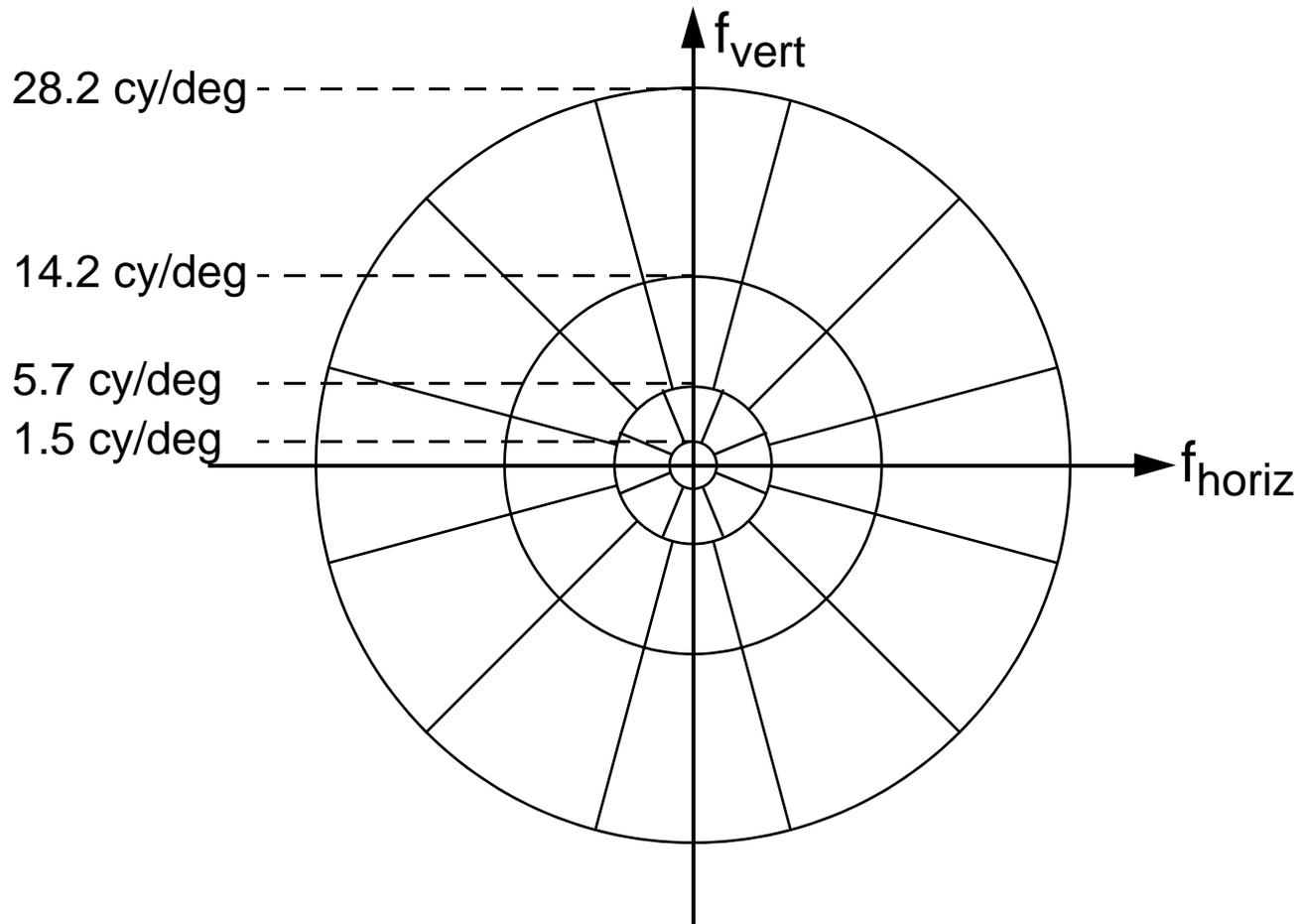
# 2-D Wavelet Transform



Input Image ($LL_0$)

Low-pass → ($LL_1$)

High-pass → ($LH_1$) ($HL_1$) ($HH_1$)

Low-pass → ($LL_2$)

High-pass → ($LH_2$) ($HL_2$) ($HH_2$)

Low-pass → ($LL_3$)

High-pass → ($LH_3$) ($HL_3$) ($HH_3$)

# Wavelet Decomposition in Frequency Space

The HVS consists of *channels*, each tuned to range of spatial frequencies and orientations.

# The Digital Signal vs. What We See

- Pixel values vs. display luminance

displayed luminance

$$L = (b + k \times p)^{\gamma}$$

monitor's luminance-to-voltage response curve exponent

black-level offset

voltage-to-pixel scaling factor

pixel value

$$Contrast = \frac{luminance\ change}{mean\ background\ luminance}$$

- We describe stimuli in terms of *contrast*.

- For complex images, we'll use *RMS contrast*.

# Contrast for Complex Images

- *RMS Contrast* defined using RMS deviation from mean background luminance $L$

$$C_{rms} = \frac{1}{\bar{L}}\sqrt{\frac{1}{N}\sum(L_i - \bar{L})^2}$$

- Recall $L = (b + k \times p)^\gamma$. For typical values of $b$, $k$, $\gamma$, this can be linearized via a Taylor series, and

$$C_{rms}^2 = \xi^2 D$$

where $D$ is the variance of the pixels, and

$$\xi = \frac{L}{k\gamma}(b + k\bar{p})^{1-\gamma}$$

# Contrast of Distorted Images



original          quantized          quantization noise

For the quantization noise, $contrast = \dfrac{1}{L}\sqrt{\dfrac{1}{N}\sum L_i^2}$

Note that we achieve $C_{max}$ at band discard.

quantized          original          quantization noise

Detection:                                                    stimuli = target

Masked detection:
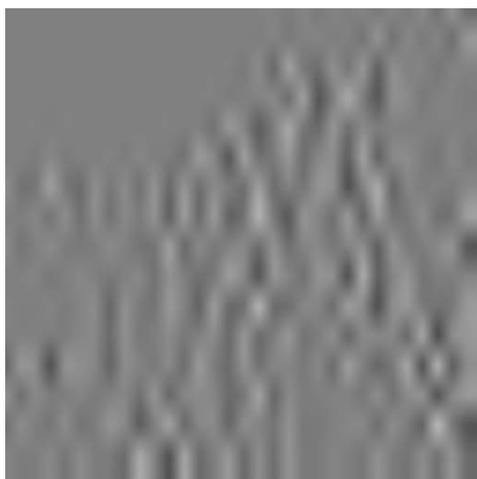
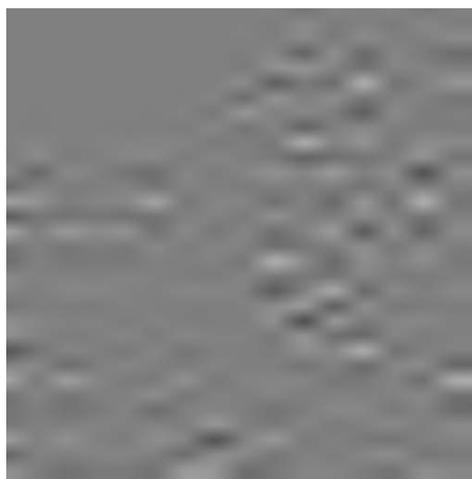stimuli          =          mask          +          target

Unmasked uniform quantization noise in the HL5, LH5, and HL5 + LH5 subbands.



<div>

If this stimulus has contrast threshold

$CT_{HL5}$

...and this stimulus has contrast threshold

$CT_{LH5}$

Then what is the contrast threshold of this stimulus?

$CT_{HL5+LH5} = \ ?$

</div>

Masked uniform quantization noise in the HL5, LH5, and  HL5 + LH5 subbands.



$$CT_{HL5} \qquad CT_{LH5} \qquad CT_{HL5+LH5} = ?$$

# Summation in Natural Images

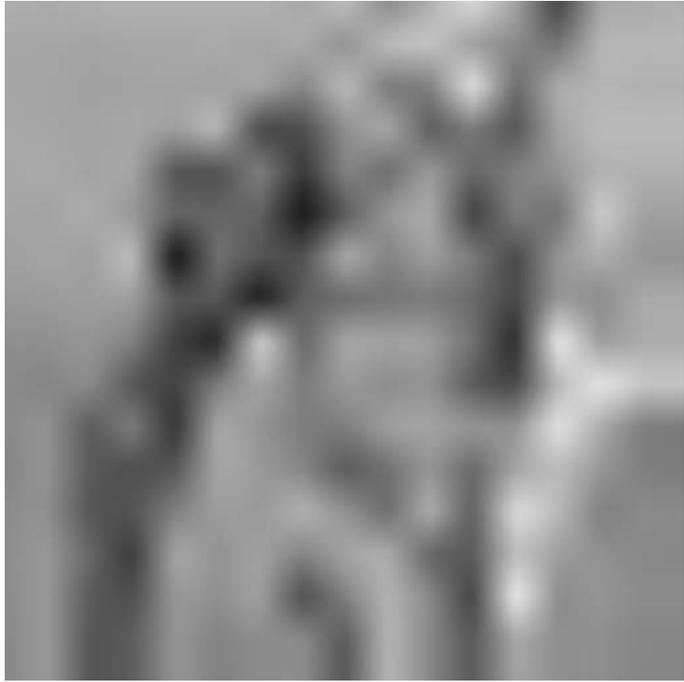- For 2 subbands simultaneously quantized in an image, $1.5 < \beta < 1.8$. Let's approximate $\beta \approx 1$.

$$(C_A / CT_A)^\beta + (C_B / CT_B)^\beta = 1$$

- *Linear* summation is consistent with summation observed in "object recognition tasks." (We are moving toward cognition...)

- *This suggests that observation is content-based rather than purely target-based* — and leads us to global precedence.
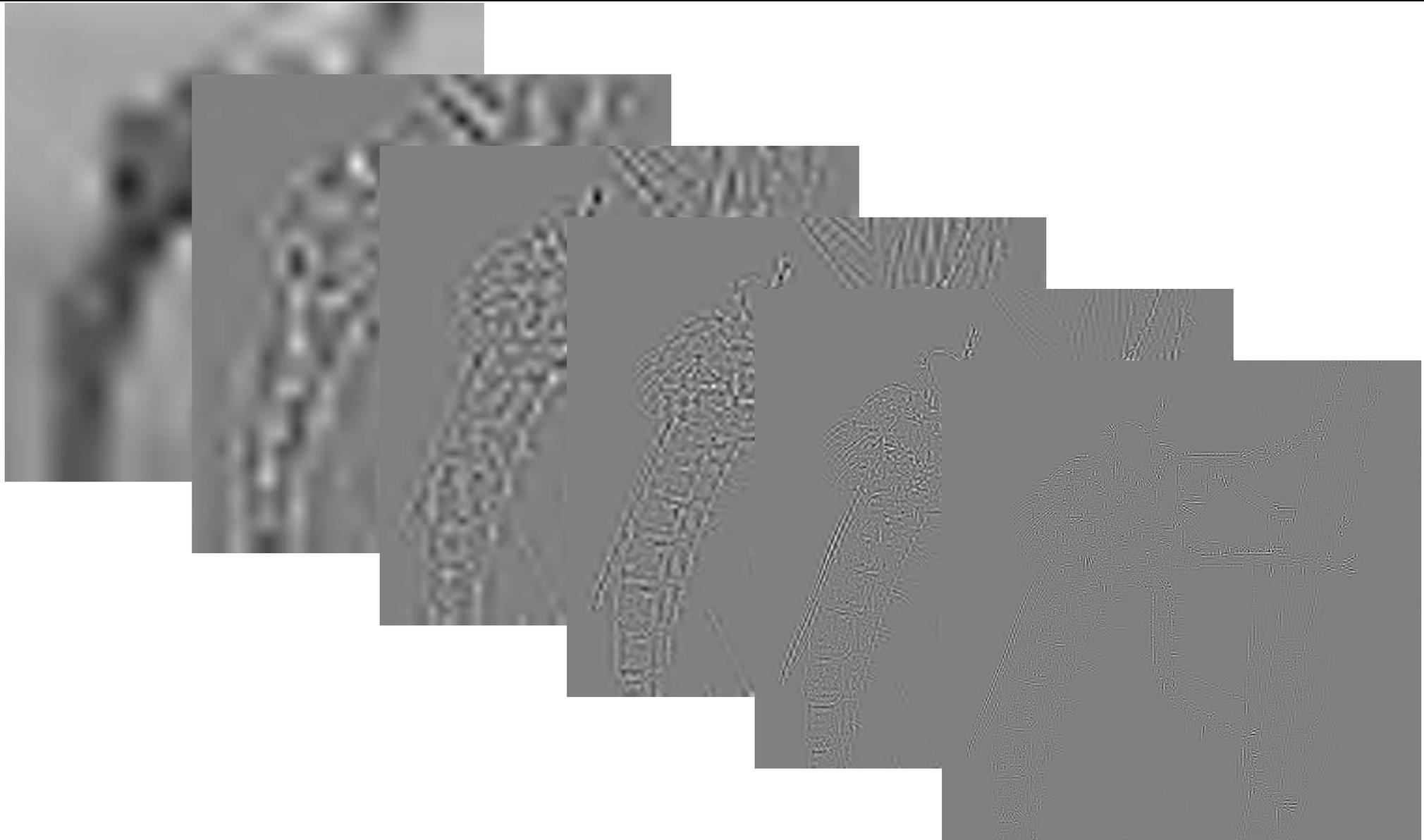
Omission of this information...

...makes us perceive this information as noise.
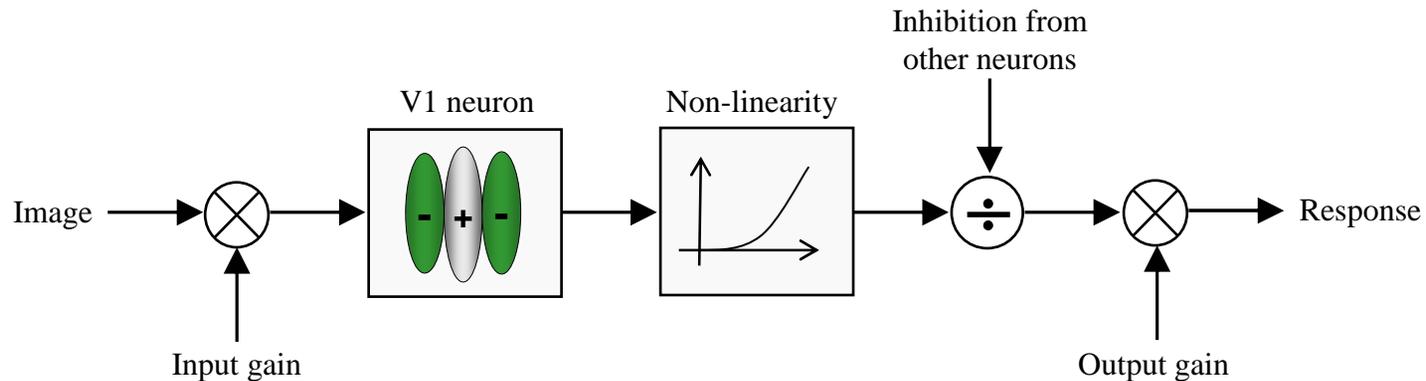
The *addition* of high-frequency content *visually degrades* the image.

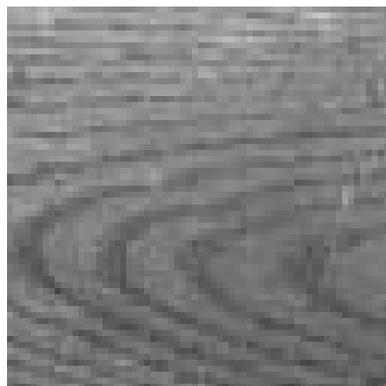# Standard Gain Control Model (Masking)



Neural response:

$$r(x, f, \theta) = \frac{w(x, f, \theta)^p}{b^q + \sum_{(x, f, \theta) \in S} w(x, f, \theta)^q}$$

$w$ = e.g., wavelet coefficient at location $x$, frequency $f$, orientation $\theta$

Usually $p \approx 2$, $q \approx 2$ — effectively variance!

Texture                    Structure                    Edge

| | | | | |
|---|---|---|---|---|
| Textures: | *fur* | *wood* | *newspaper* | *basket* |
| Structures: | *baby* | *pumpkin* | *hand* | *cat* *flower* |
| Edges: | *butterfly* | *sail* | *post* | *handle* *leaf* |

- Experimentally quantify masking of texture/structure/edge patches, and develop an appropriate gain control model.

*Textures mask more than structures (2x), which mask more than edges (2.5x).*

Neural response:

$$r(x, f, \theta) = \frac{w(x, f, \theta)^p}{b^q + g_m \sum_{(x, f, \theta) \in S} w(x, f, \theta)^q}$$

$g_m$ is an inhibitory modulation term and varies based on patch type

Texture                Structure                Edge

- Masked CSF and summation/global precedence

  - Distortion-contrast quantization.

  - A new multiple description quantization strategy.

  - Visual signal-to-noise ratio (VSNR) — a quality metric.

- Masked CSF, summation/global precedence, and gain control model

  - Overhead-free optimal spatially localized quantization.

# Distortion-Contrast Quantization

A quantization strategy for wavelet-coded natural images based on

1. Our masked detection results at and above threshold;

2. Linearity in summation;

3. Global precedence.

Result: a strategy which works seamlessly for all rates, producing better looking images at up to 30% lower rates.

JPEG-2000 compatible (but not necessary!)

# Harbor, 0.4 bits/pixel, JPEG-2000 Framework
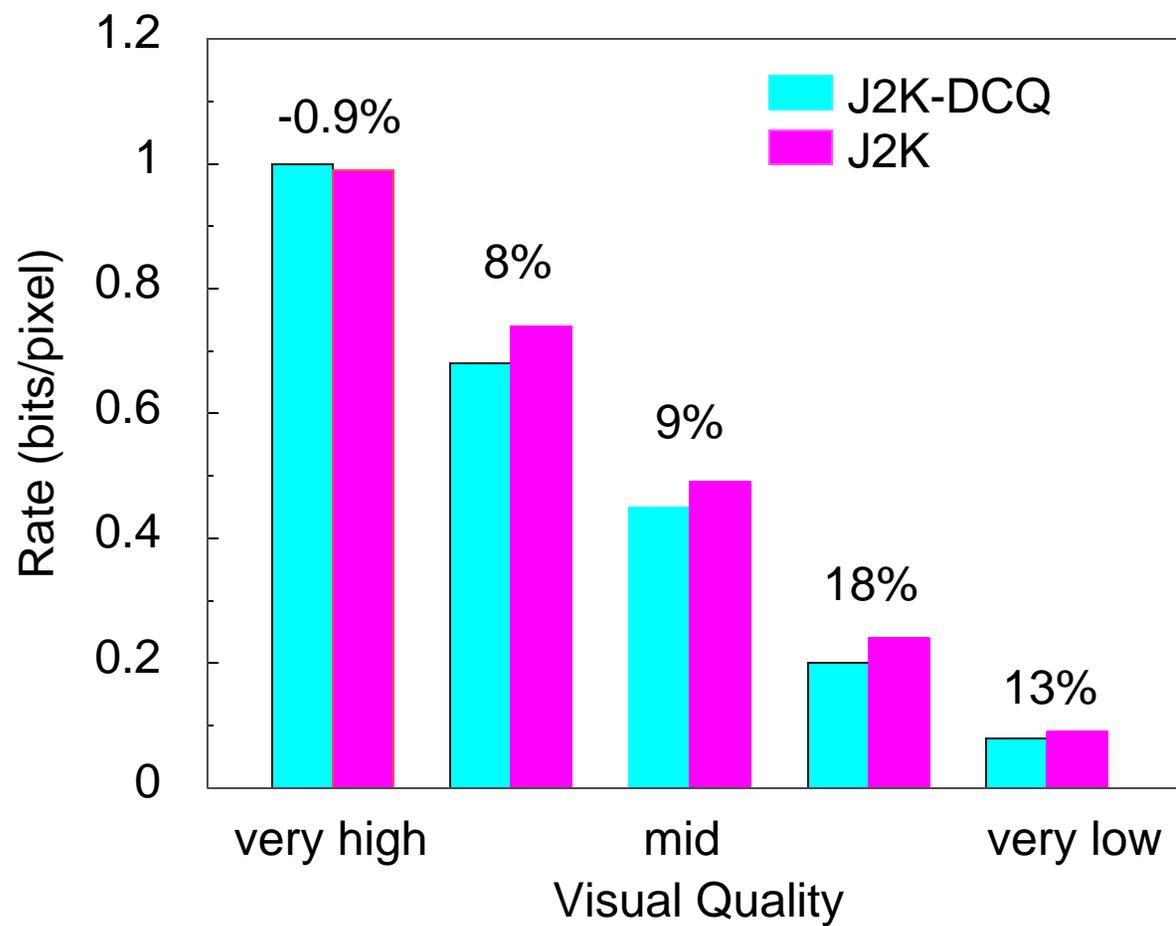
JPEG-2000

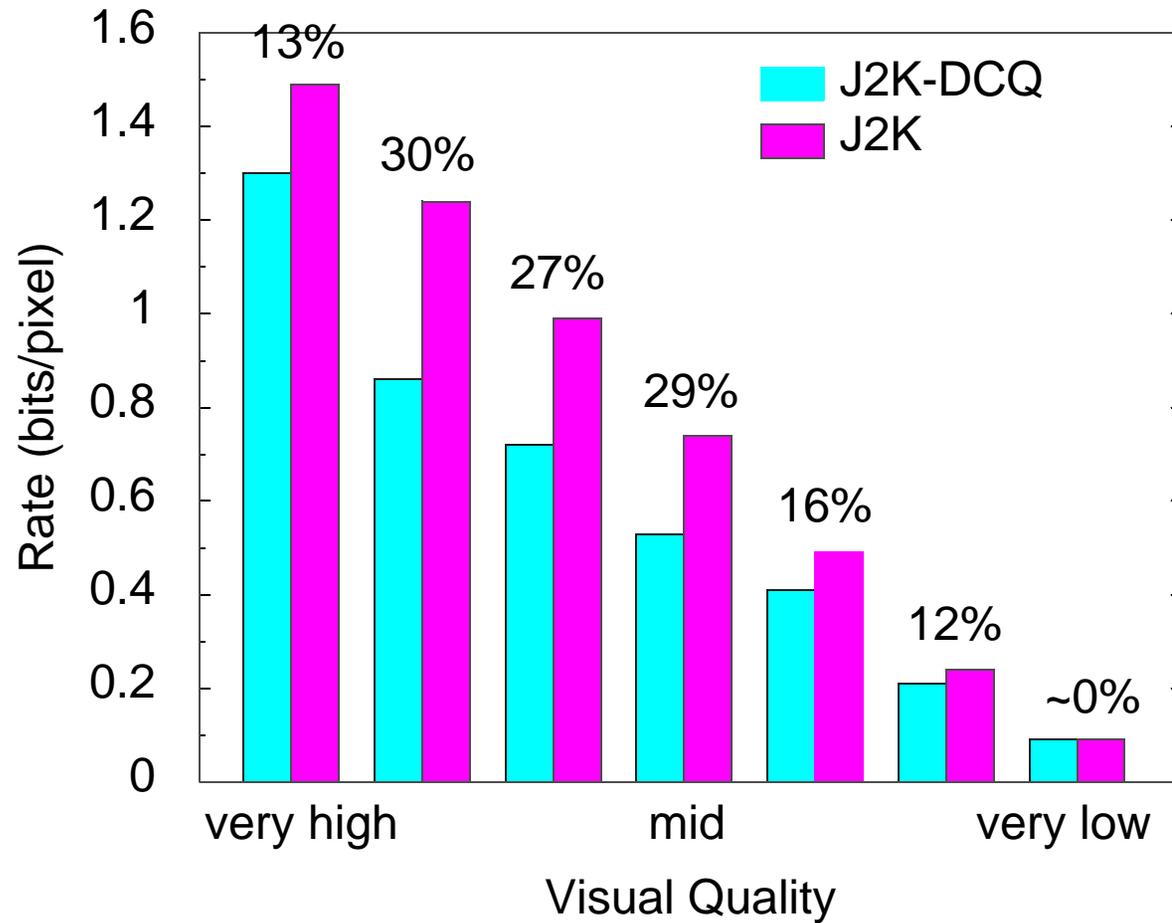Contrast-based JPEG-2000

Cat

Rainriver

# At Equal Quality: Rate Savings for Cat

# At Equal Quality: Rate Savings for Rainriver

# Overhead-free Optimal Spatially Localized Quantization

- Goal — set quantization step sizes locally within an image according to local masking thresholds.

- Problem — step sizes must then be transmitted along with the image. Until now, the overhead has proved to be prohibitive.

- Our solution — information used to produce the step sizes is used as side information to compress the image. This does NOT incur a rate penalty: conditioning reduces entropy.

# Visual & PSNR Results

| Distortion visibility | Image | Number preferred | | PSNR | |
|---|---|---|---|---|---|
| | | Proposed | JPEG-2K | Proposed | JPEG-2K |
| Barely visible | horse (1.13 bpp) | 8 | 0 | 30.0 | 32.1 |
| | rhino (1.88 bpp) | 7 | 1 | 24.3 | 29.0 |
| Very visible | horse (0.64 bpp) | 6 | 2 | 27.0 | 28.3 |
| | rhino (1.24 bpp) | 6 | 2 | 21.3 | 25.7 |

Spatially localized quantization hides much more error in the image for the same visual quality.

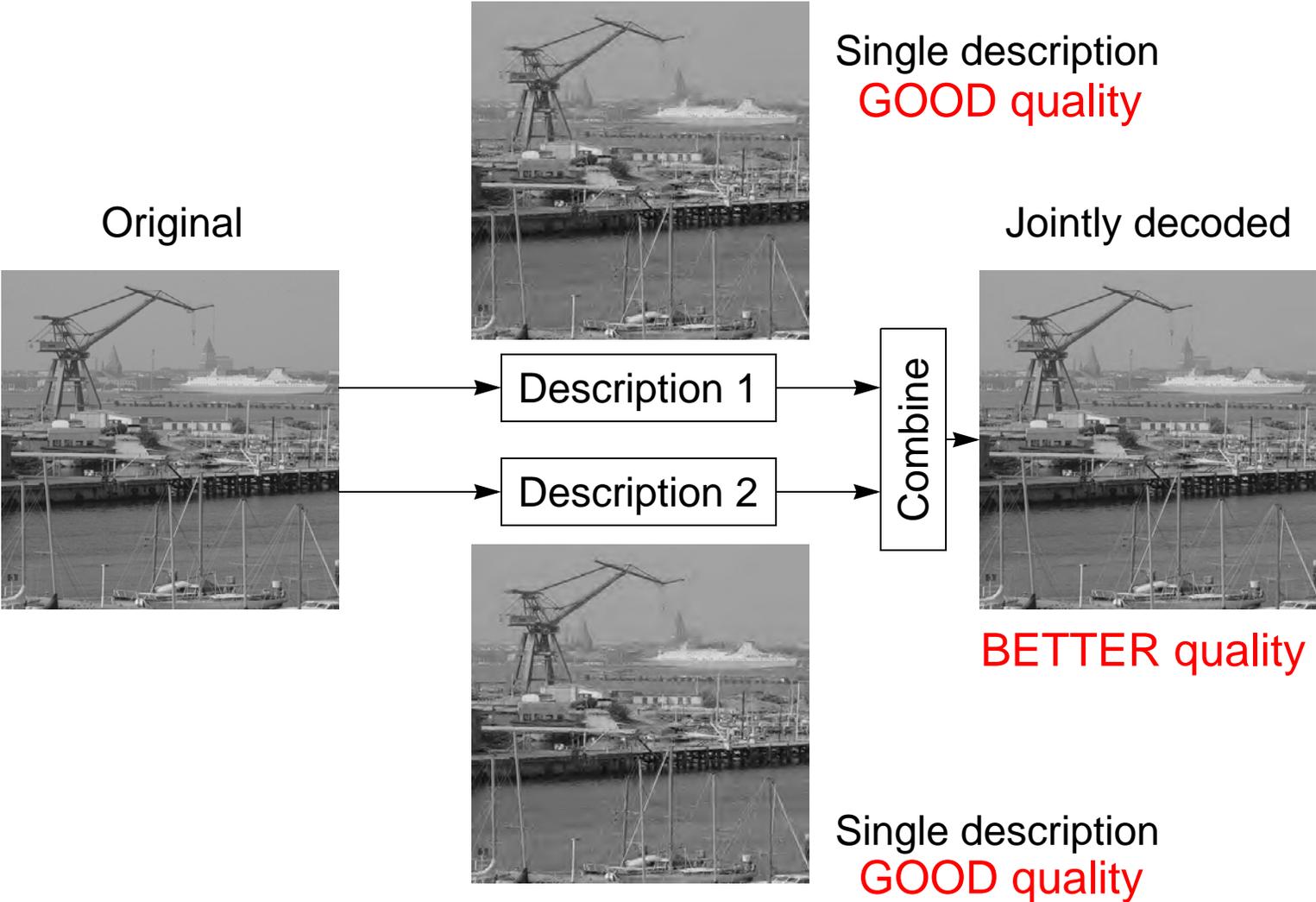# Multiple Description Image Coding



Original

Single description
GOOD quality

Jointly decoded
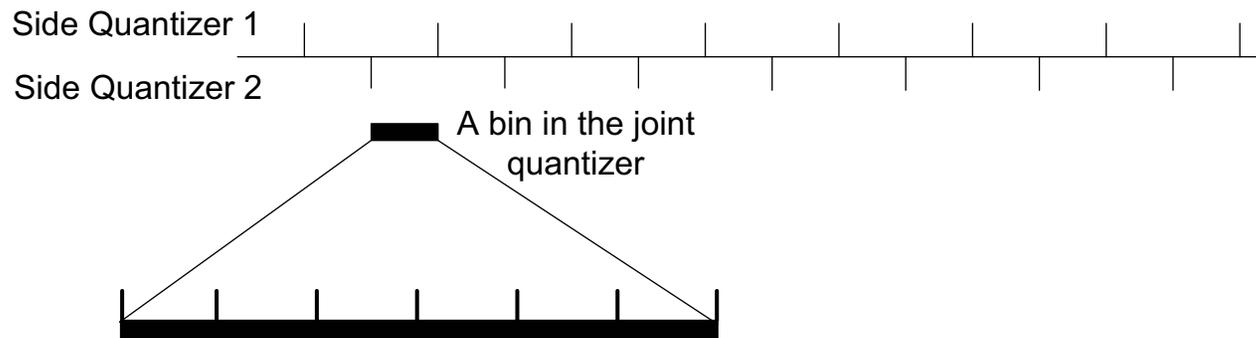
Description 1

Description 2

Combine

BETTER quality

Single description
GOOD quality

# Visually Optimized MD Image Coding

- Problem: HVS results are for distortions caused by uniform (convex) quantization cells, BUT "standard" MD quantizers use non-convex cells.

- Our solution: design a new MD quantization strategy which has equivalent R-D performance to standard techniques but which uses convex cells.

Side Quantizer 1

Side Quantizer 2

A bin in the joint quantizer

MSE-optimized

Visually-optimized

MSE-optimized                    Visually-optimized

# Concluding Comments

- Extensive psychophysical experiments have yielded more accurate HVS characterizations for image compression.

- These HVS characteristics have been used to drive signal processing algorithm development.

- The resulting algorithms outperform current state-of-the-art results.

- We have also applied this methodology to the design of a video quality measure.

*What is task-based imaging?*

From the user/application perspective:

- Who is viewing it and why?

- How is the visual information to be used?

From the image processor's perspective:

- What *must* be conveyed by the visual information?

- What is nice to have, but optional for the task?

quantized = original + quantization noise

Detection: stimuli = target

Masked detection:
stimuli = mask + target

**The target is the *distortion***

= quantized    original    + quantization noise

**Detection:**    stimuli = target

**Masked detection:**
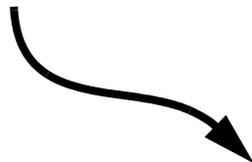stimuli    =    mask    +    target

# Questions

- How can we measure "usefulness" of an image?

- What distortions should we explore?

- How is "quality" related to usefulness (utility)?

- Can current quality estimators predict utility?
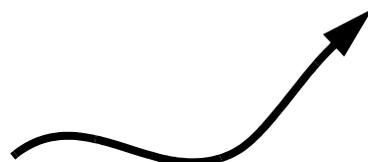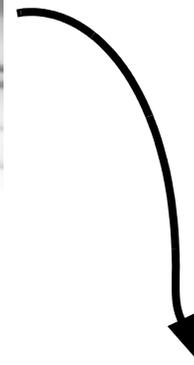
- Can we create a *utility estimator*?

Not recognizable

Recognizable
but distorted

Recognition
threshold

Visually
lossless

Not recognizable

Recognizable
but distorted

Recognition
threshold

Visually
lossless

Not recognizable

Recognizable
but distorted

Recognition
threshold

Visually
lossless

How "long" are these distances?

1. Recognition

   Single-image stimulus: "Do you recognize the image content?"

2. Utility assessment

   Image pair stimulus: "Which image tells you more about the content?"

3. Quality assessment

   SAMVIQ or ACR

# Graduate Student Acknowledgements

- HVS/Perception: Prof. Damon Chandler, Dr. Marcia Ramos, Dr. Mark Masry, Bobbie Chern, Jeri Moses.

- Image Applications: Prof. Damon Chandler, Dr. Matt Gaubatz, Dr. Chao Tian.

- Utility work: Dr. David Rouse.

Papers on all these topics can be found at

http://foulard.ece.cornell.edu