# Auditory Attention: From Saliency to Models (and Applications)

Malcolm Slaney
June 7, 2017

# Binaural Workshop

# The Problem

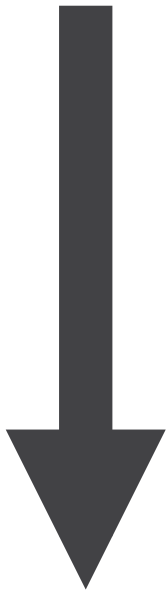- **A Cocktail Party**

- **The Open-Microphone Problem**

# Types of Attention

- **Top down**
  Driven by needs,
    experience
  Language models

- **Bottom Up**
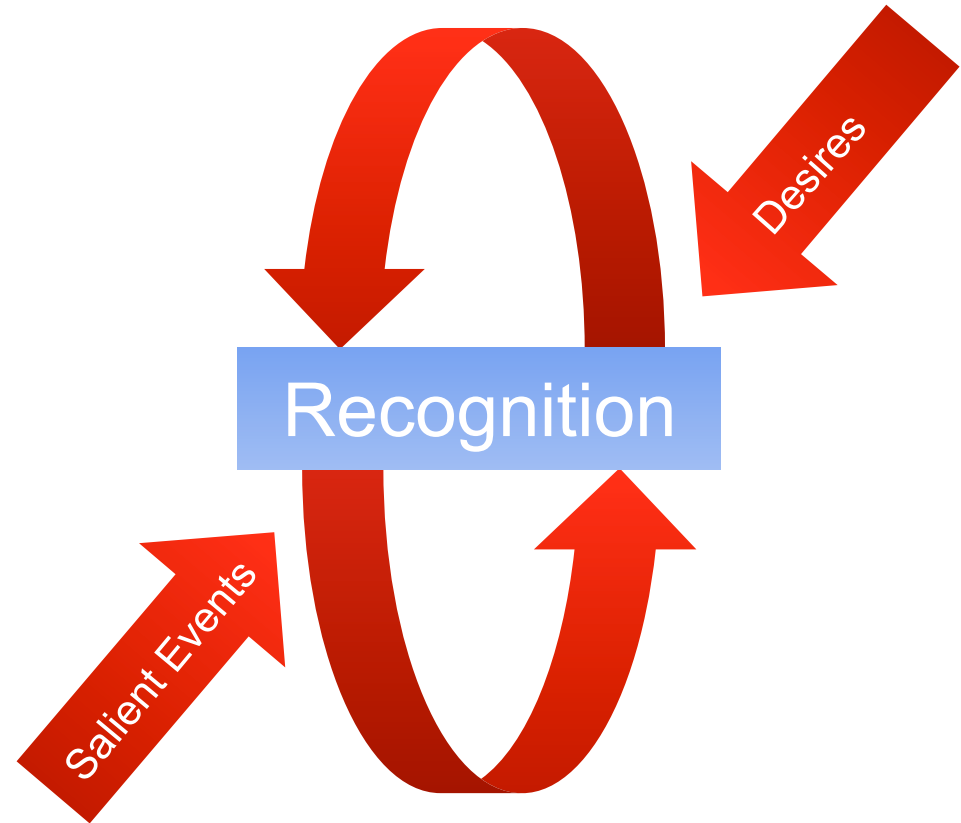  Driven by perceptual
    surprise
  Saliency

Simplification

# Outline

**Introduction**
- Top Down
- Breaking the loop
- Eyes vs. ears

**Saliency**
- Data
- Models

**Attention**
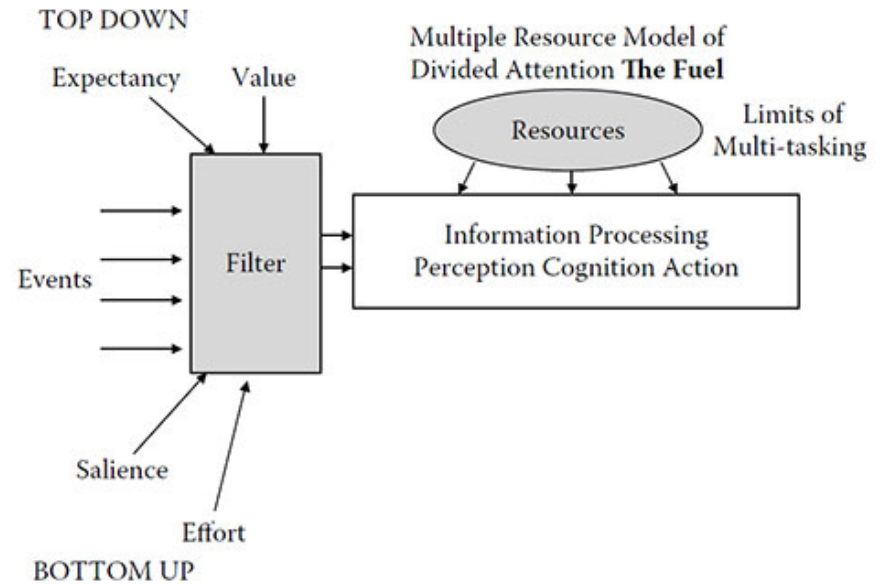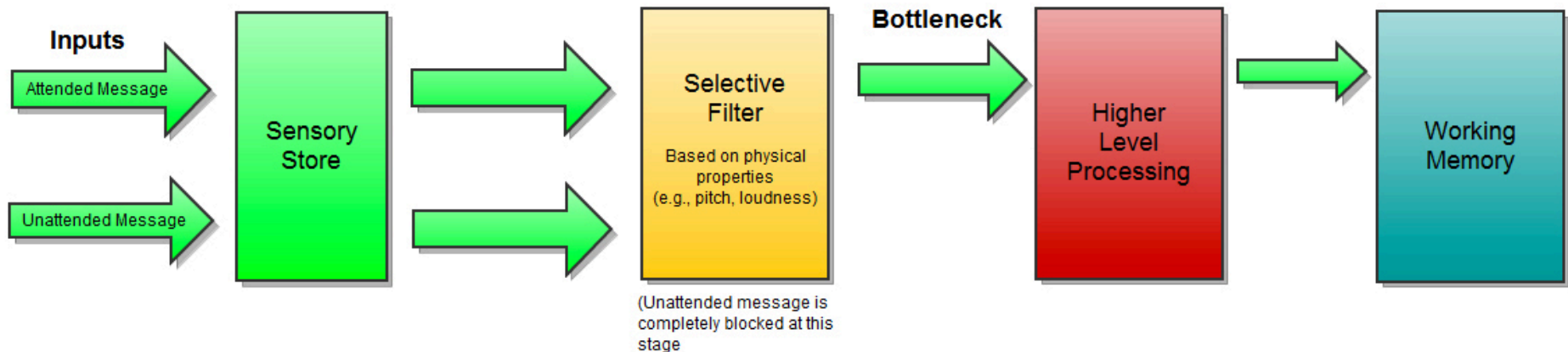- Models
- Decoding
- Applications

# Role of Attention?

**Limits**
- Sensory filter
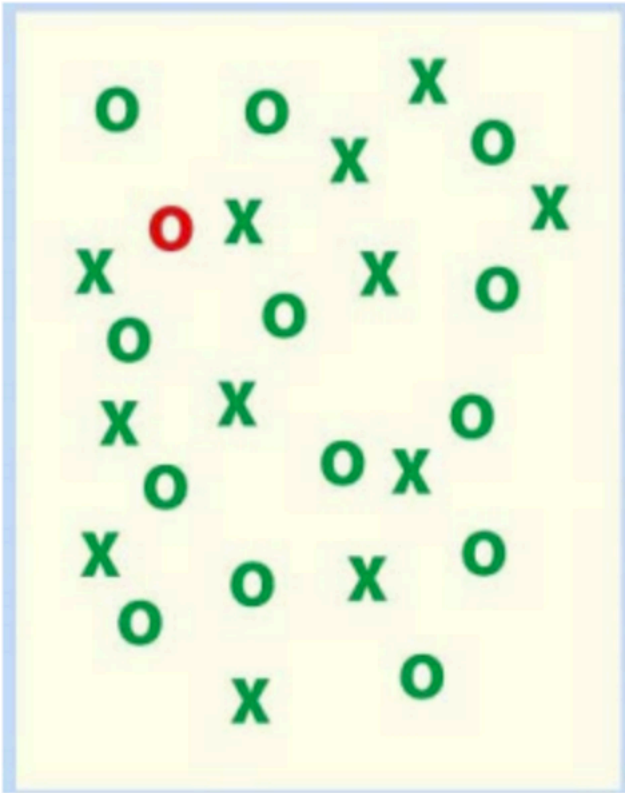- Cognitive resources

**Integrative mechanism**
- Bind features

TOP DOWN

Expectancy    Value

Events    →    Filter    →    Information Processing
Perception Cognition Action

Multiple Resource Model of
Divided Attention **The Fuel**

Resources    Limits of
Multi-tasking

Salience

Effort

BOTTOM UP

*Broadbent's Filter Model*

**Inputs**

Attended Message

Unattended Message

Sensory
Store

Selective
Filter

Based on physical
properties
(e.g., pitch, loudness)

(Unattended message is
completely blocked at this
stage)

**Bottleneck**

Higher
Level
Processing

Working
Memory

# Popout



Bottom Up
(Popout)

Top Down
(Search)

# Eye Tracking

Input image

Multiscale low-level feature extraction

Colours
Red, green, blue, yellow, etc.

Intensity
On, off, etc.

Orientations
0°, 45°, 90°, 135°, etc.

Other
Motion, junctions and terminators, stereo disparity, shape from shading, etc.

Centre–surround differences and spatial competition

Feature maps

Feature combinations

Top-down attentional bias and training

Attended location

Inhibition of return

Winner-take-all

Saliency map

Itti & Koch, Nature Reviews Neurosci 2001
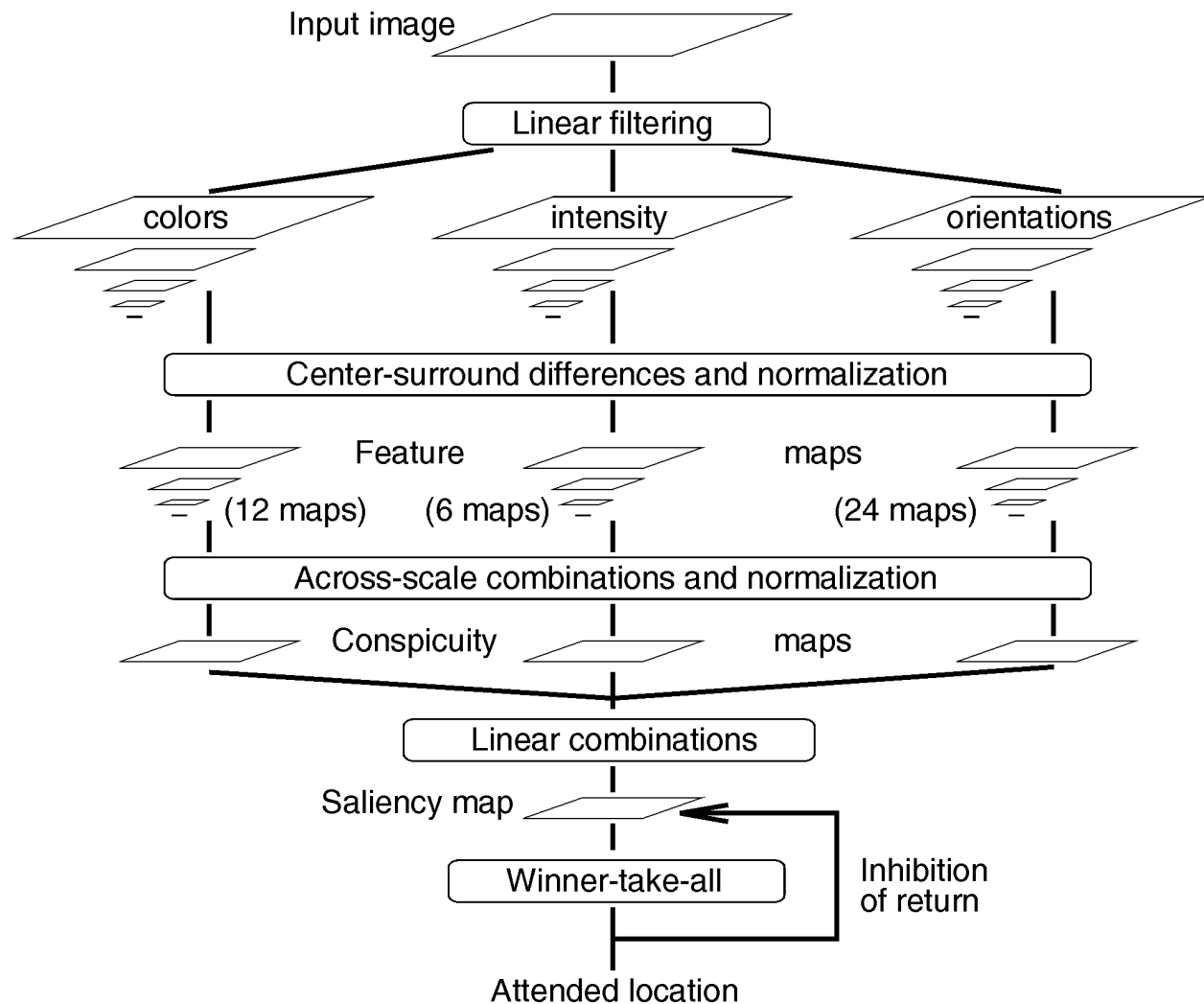
# Itti – Salience Model

**Combines**
- Color
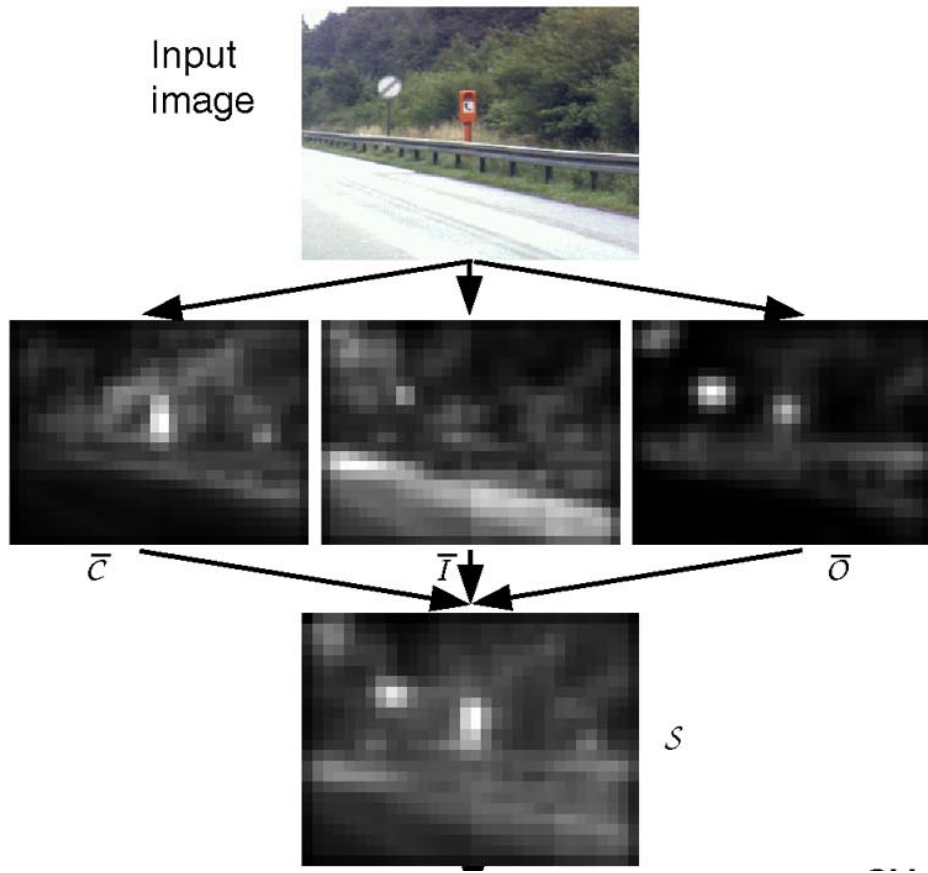- Intensity
- Orientation

**Processing**
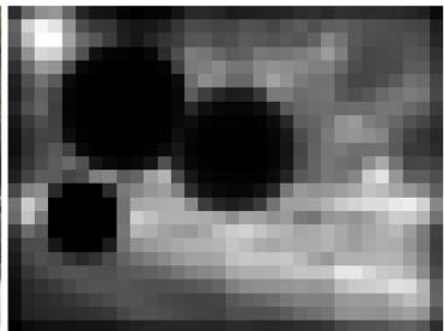- Center-surround differences
- Normalization
- Multi-scale

**Decision**
- Linear combination
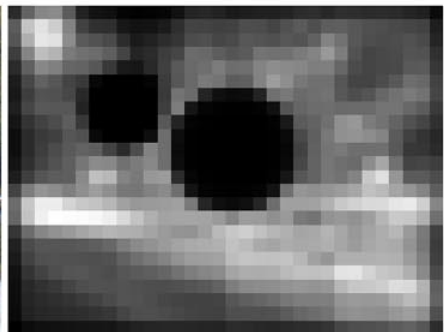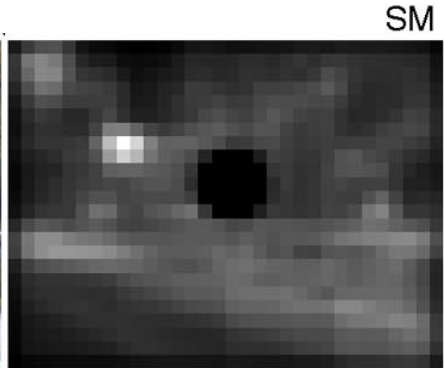- Winner take all
- Inhibit and repeat

Input image

Linear filtering

colors    intensity    orientations

Center-surround differences and normalization

Feature           maps

(12 maps)    (6 maps)           (24 maps)

Across-scale combinations and normalization

Conspicuity           maps

Linear combinations

Saliency map

Winner-take-all

Inhibition of return

Attended location

# Detecting Attention

- **Eye Gaze**

- **Ear Gaze**





>24 databases
  Images vs. Movies

VS

# Salient Sounds

# Focusing on the clutter in auditory scenes: Perspectives from modeling auditory attention

### Emine Merve Kaya and Mounya Elhilali[*]

*Laboratory for Computational Audio Perception, Department of Electrical and Computer Engineering, the Johns Hopkins University, 3400 N Charles street, Barton Hall, Baltimore, MD 21218, USA, orcid.org/0000-0003-2597-738X*

# Kayser Test Sounds

Which is more salient?

# Salient Sound Detection (Elhilali at JHU)

## Musical Examples

stimA4_all

stimA4

aro15example

## Speech Examples
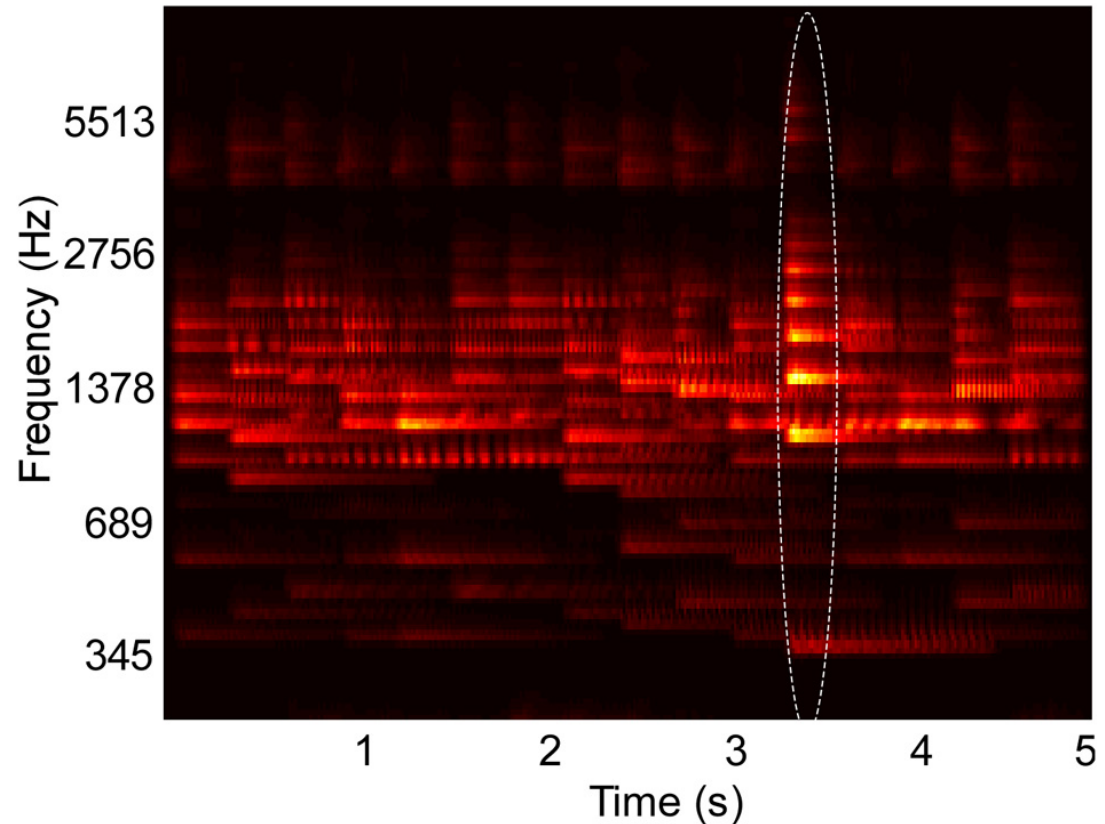
Amini_all

Amini_timbrepitch

# Kaya – Human Saliency Tests

## Background

- Overlapping musical notes
- Pitch and intensity constrained
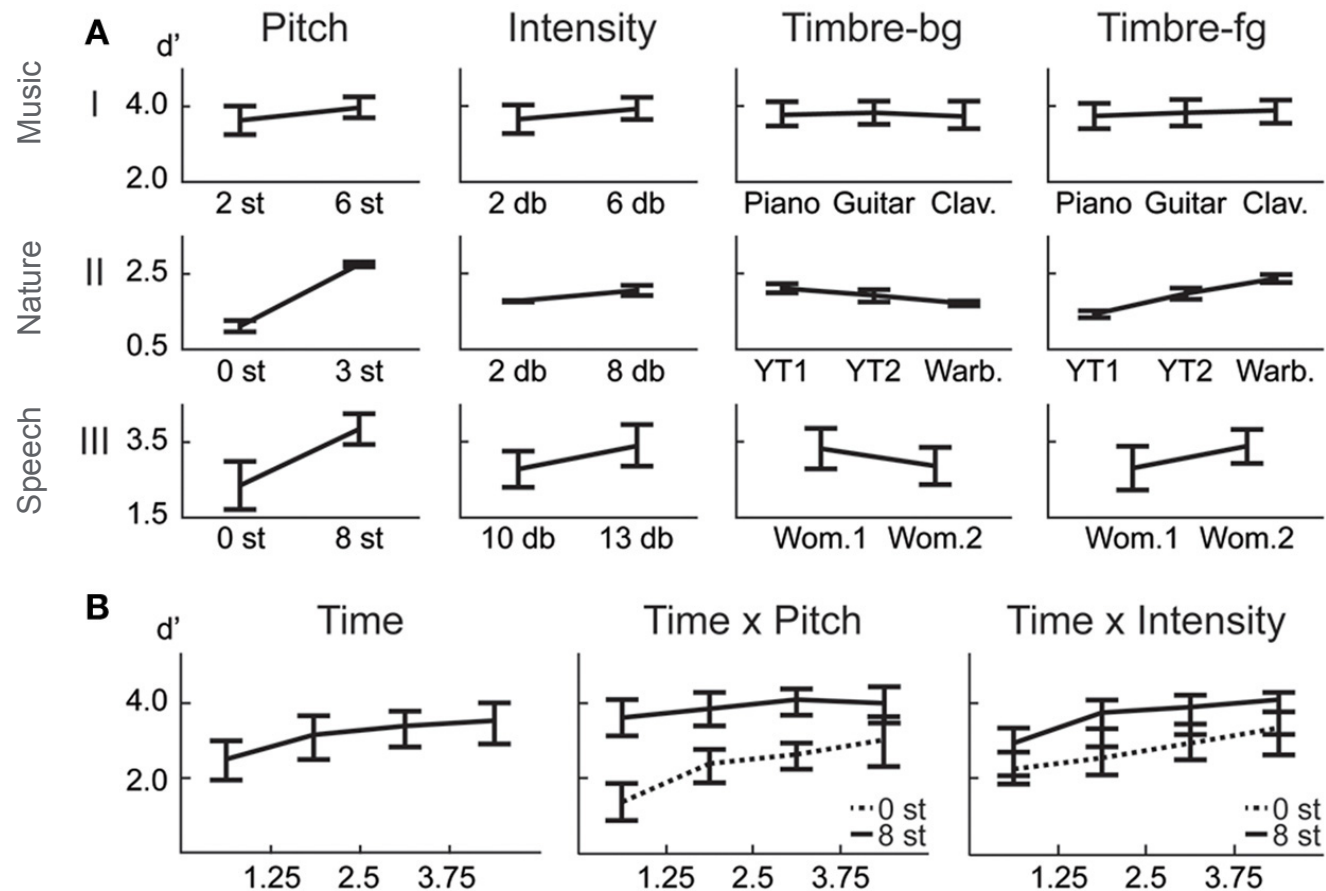- Pitch from 196-247Hz

## Foreground

- 350Hz, +6dB

# Kaya – Human Saliency Tests

**Detectability**

**d' measures separation between the means of the signal and the noise distributions, compared against the standard deviation of the signal plus noise distributions.**
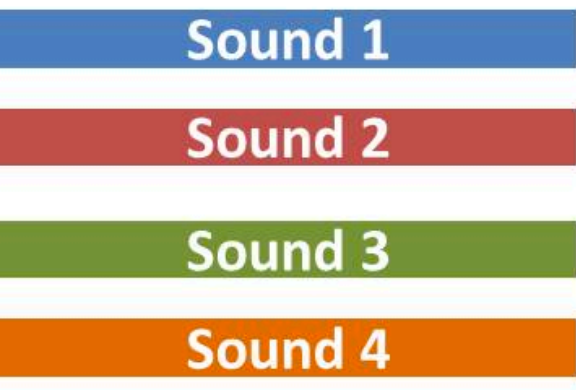
# Yahoo Captchas (Unpublished Pilot)

**Objective measure of salience**
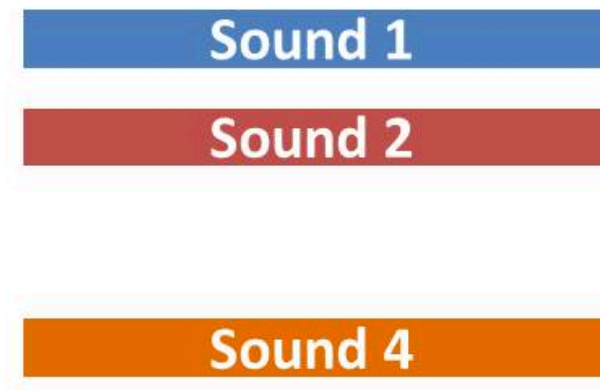
• Background speech babble
• Recognize foreground digits

# Distractors (Maria Chait at UCL)

# Question?

**Do we care more about distractors or detectability?**

**Detectability**

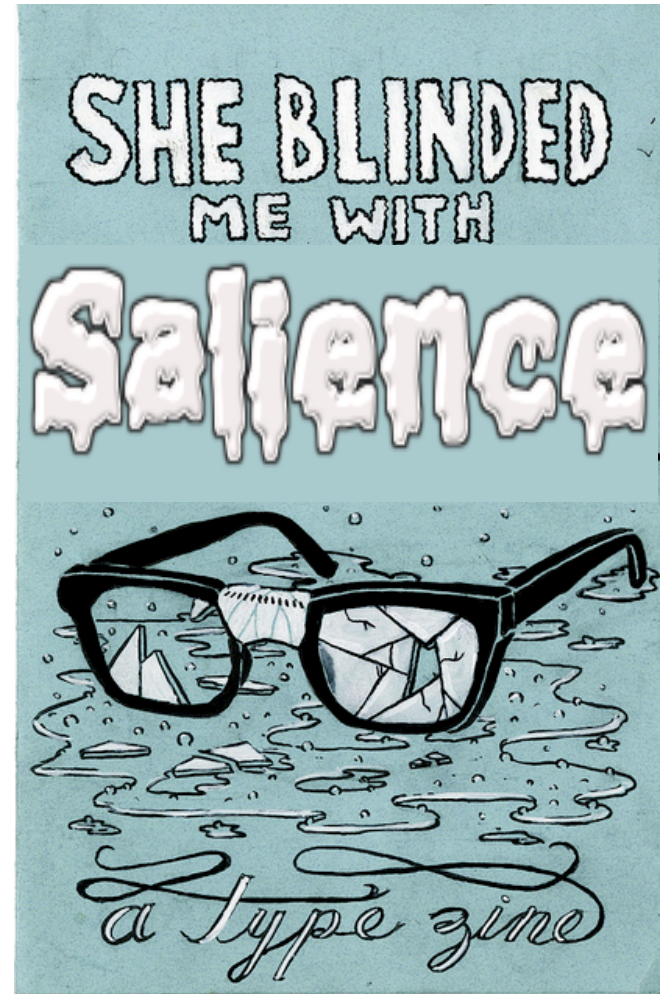- **Can we hear the difference?**

- **Precursor to distraction?**

**Distractors**

- **More ecological**

- **Did it change your attention?**
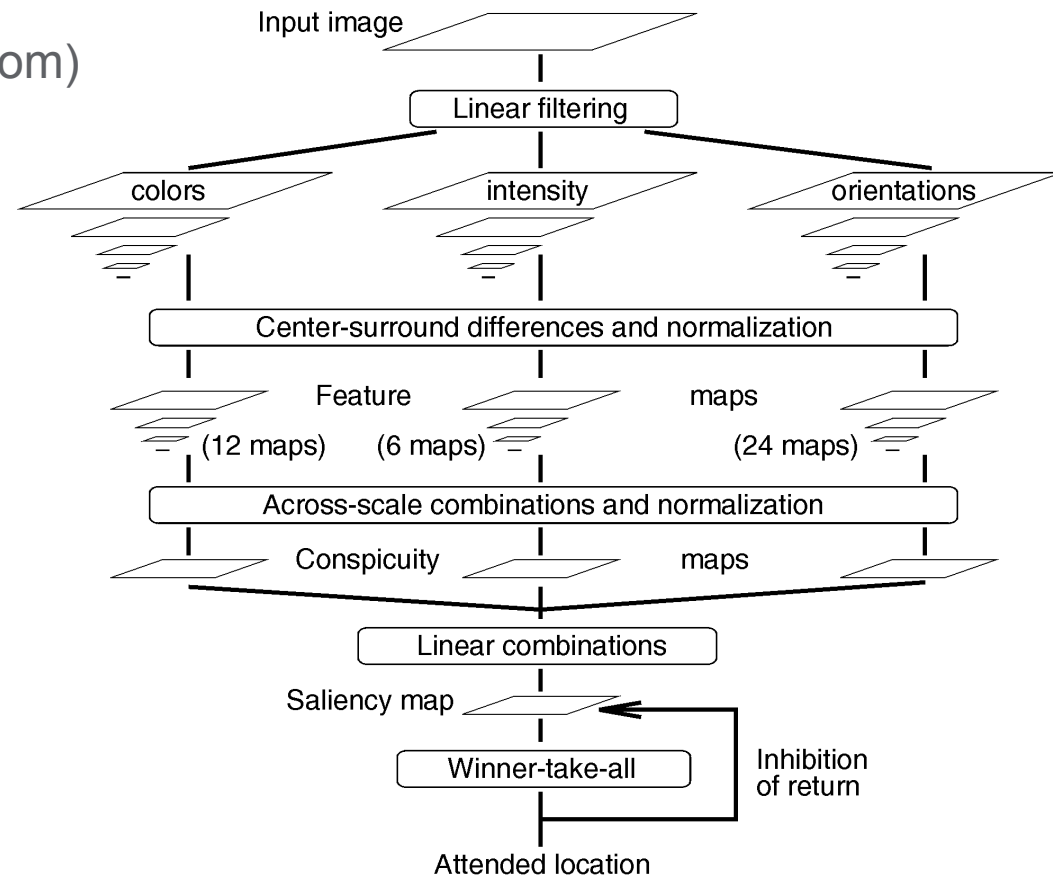
# Bottom Up Models

# Visual models

- Direct analogue (Kayser)
- Add pitch & orientation (Kalinli)
- Change to modulation (Duangudom)
- Use entropy (Wang)

# Temporal

- Add time (Kaya)
- Add tracking (Kaya)
- Use statistics (Tsuchida)

# Machine learning

- Learn from meetings (Kim)

Input image

Linear filtering

colors    intensity    orientations

Center-surround differences and normalization

Feature        maps

(12 maps)    (6 maps)        (24 maps)

Across-scale combinations and normalization

Conspicuity        maps

Linear combinations

Saliency map

Winner-take-all

Inhibition
of return

Attended location

# Kayser's Saliency Model
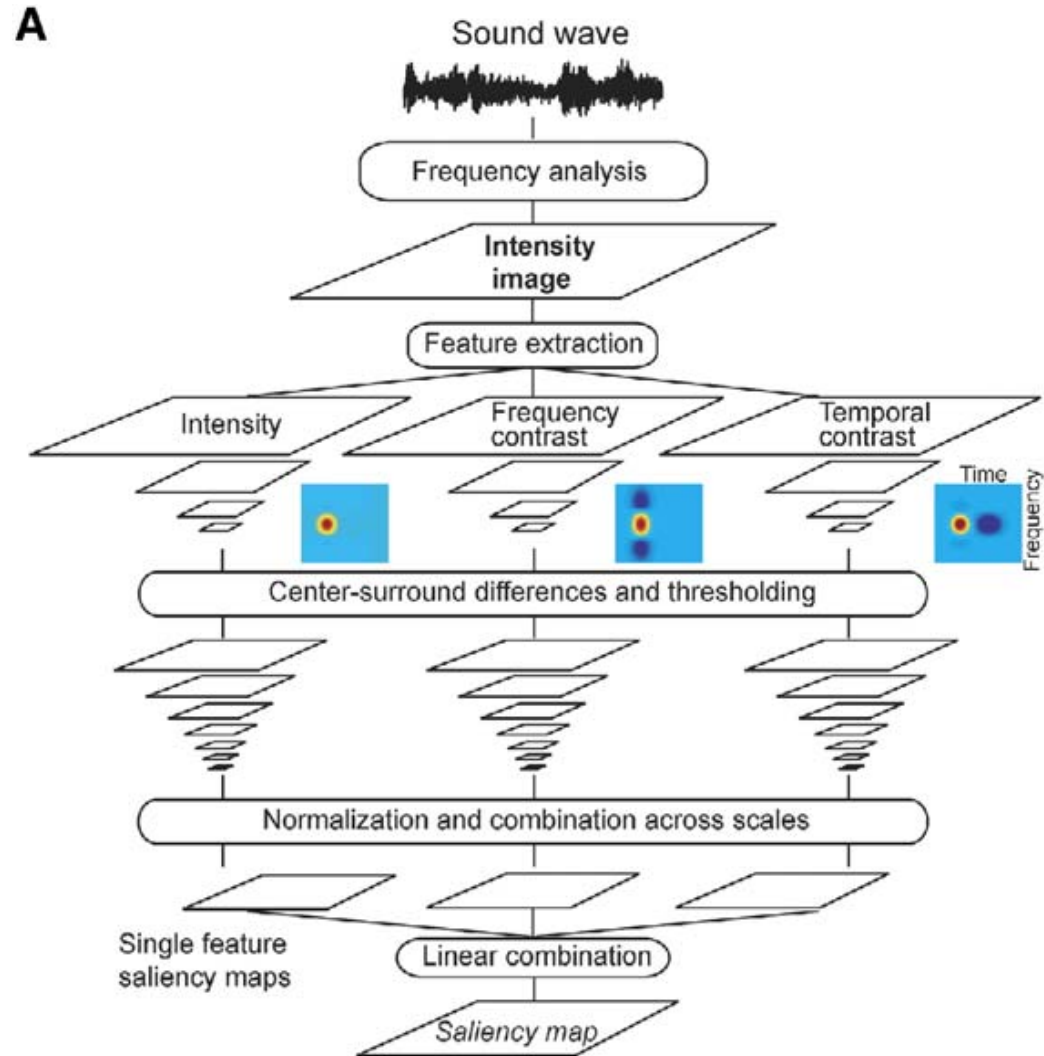
**Direct analogue to Itti Model**
- Spectrogram
- Intensity map
- Frequency contrast
- Temporal contrast

**Processing**
- Multiscale
- Center-surround differences
- Thresholding
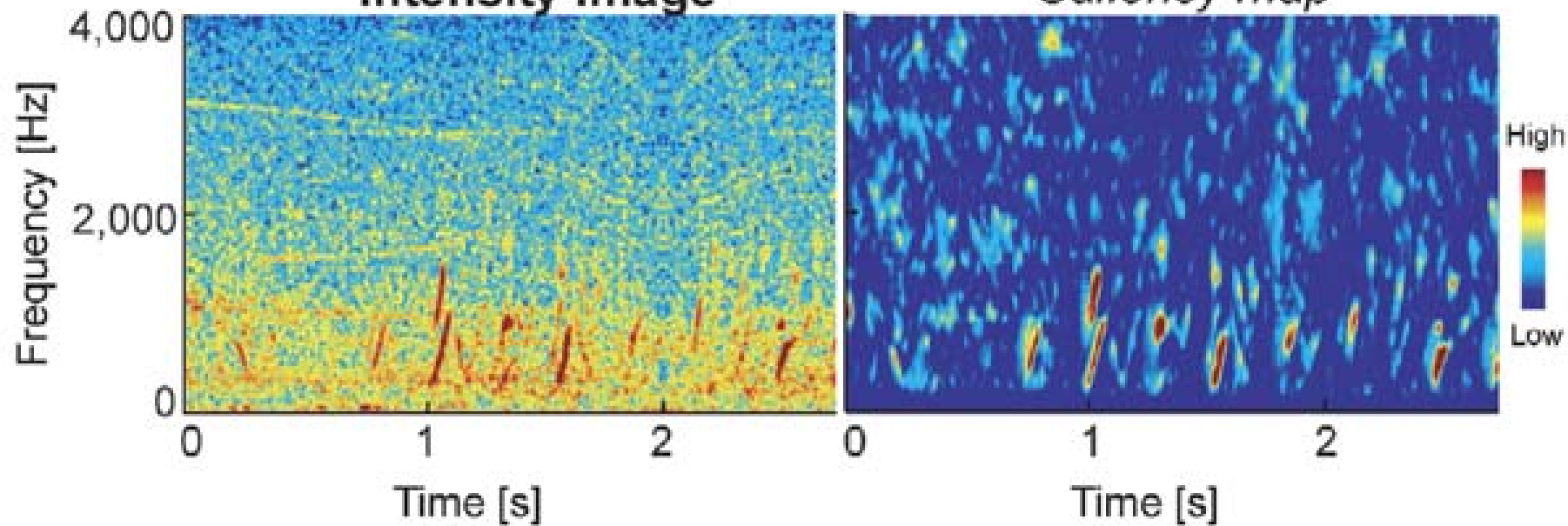
**Decision**
- Linear weights
- Form map

A

Sound wave

Frequency analysis

**Intensity image**

Feature extraction

Intensity | Frequency contrast | Temporal contrast

Time / Frequency

Center-surround differences and thresholding

Normalization and combination across scales

Single feature saliency maps

Linear combination

Saliency map

Kayser C, Petkov CI, Lippert M, Logothetis NK. Mechanisms for allocating auditory attention: an auditory saliency map. Curr Biol 2005;15(21):1943-1947.

# Kayser's Example

**Spectrogram**          **Saliency Map**



Kayser C, Petkov CI, Lippert M, Logothetis NK. Mechanisms for allocating auditory attention: an auditory saliency map. Curr Biol 2005;15(21):1943-1947.

# Kayser – More examples

**Tones**
- Salient irrespective of length
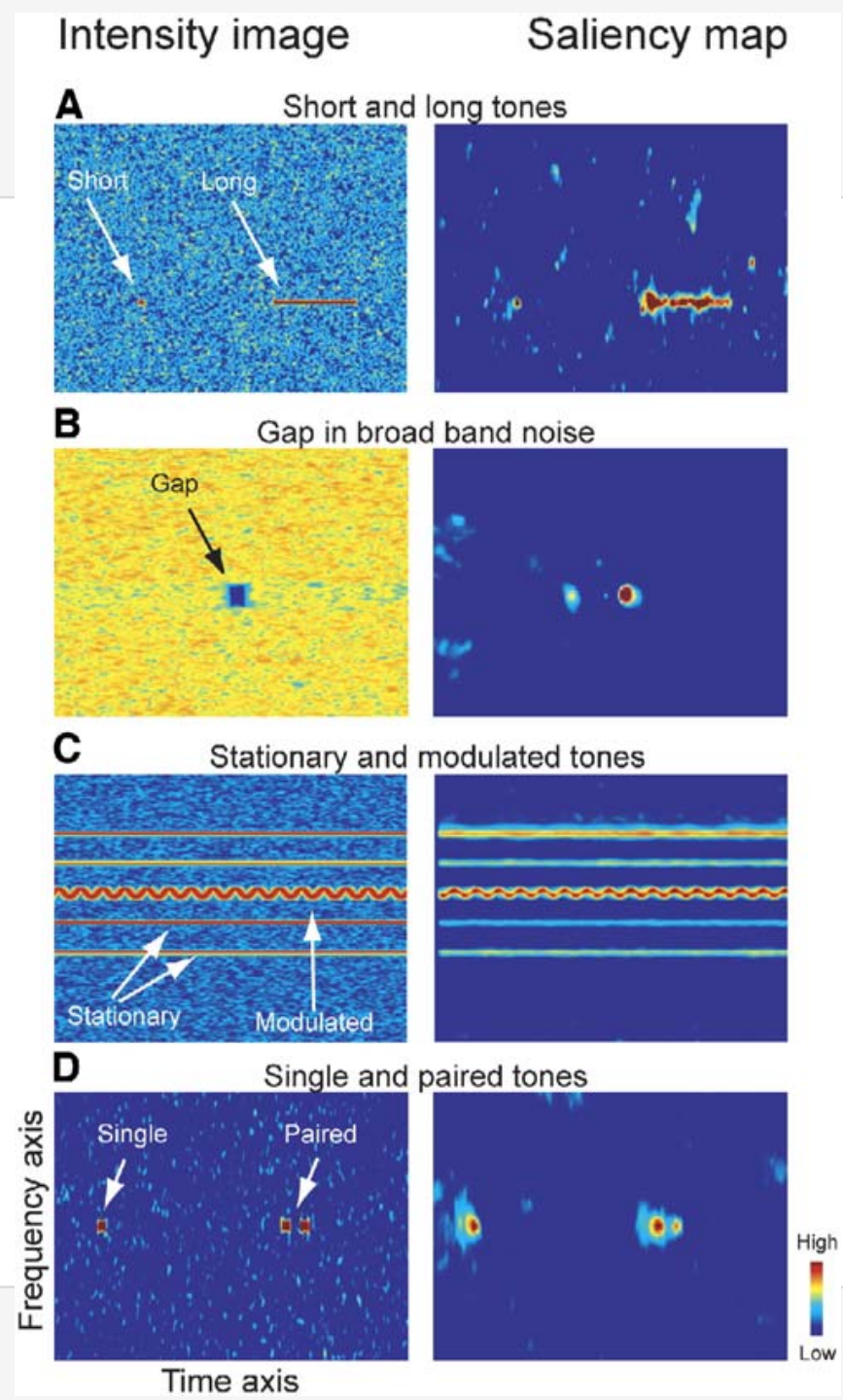- Longer events accumulate higher saliency

**Gaps**
- Missing part is salient

**Modulation**
- Modulated events are more salient than stationary
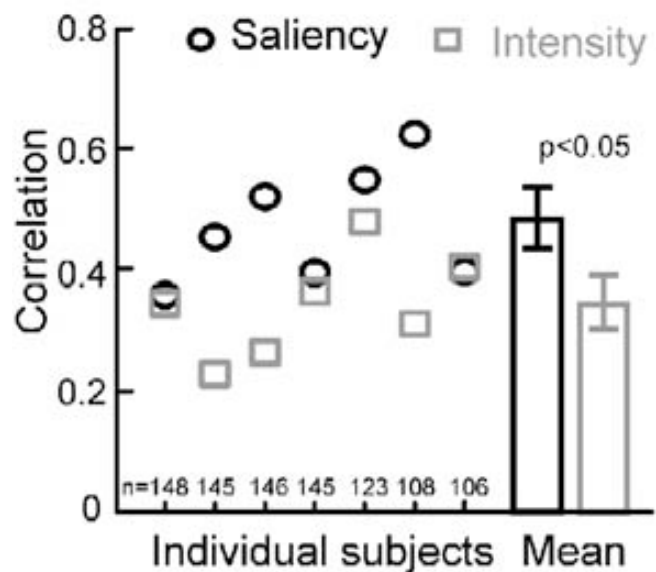
**Forward masking**
- Second is less salient

Kayser C, Petkov CI, Lippert M, Logothetis NK. Mechanisms for allocating auditory attention: an auditory saliency map. Curr Biol 2005;15(21):1943-1947, Supplement.
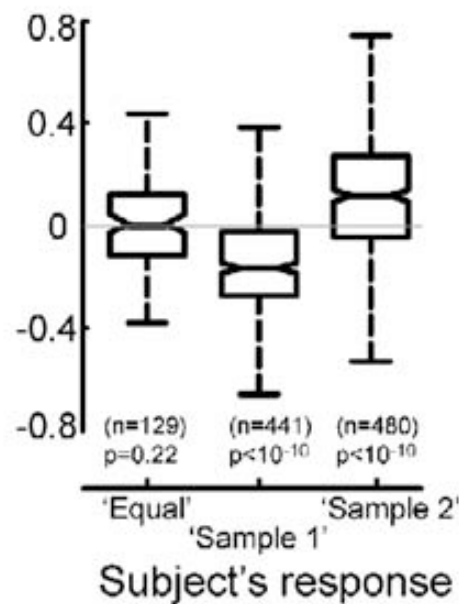
# Kayser Human Tests

## Human tests

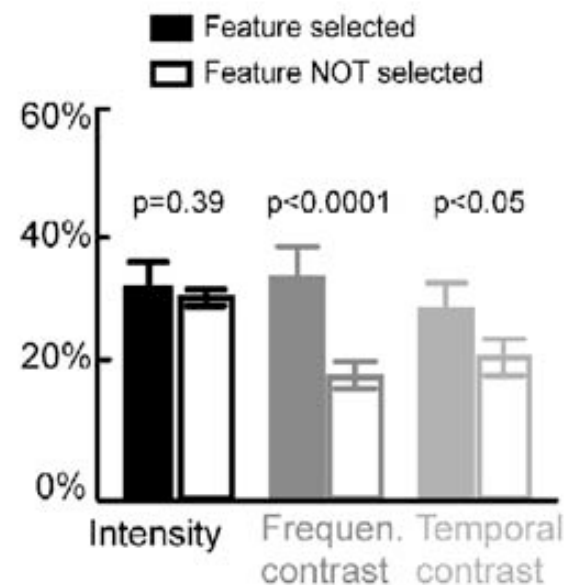- Present two examples
- Just higher saliency
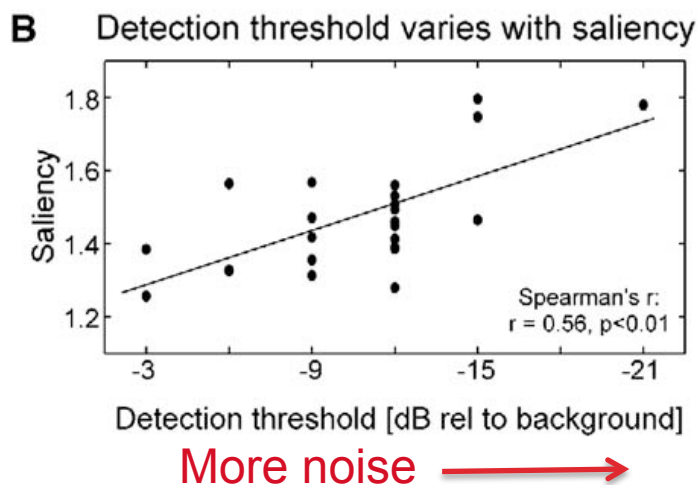


Correlation of subject & model
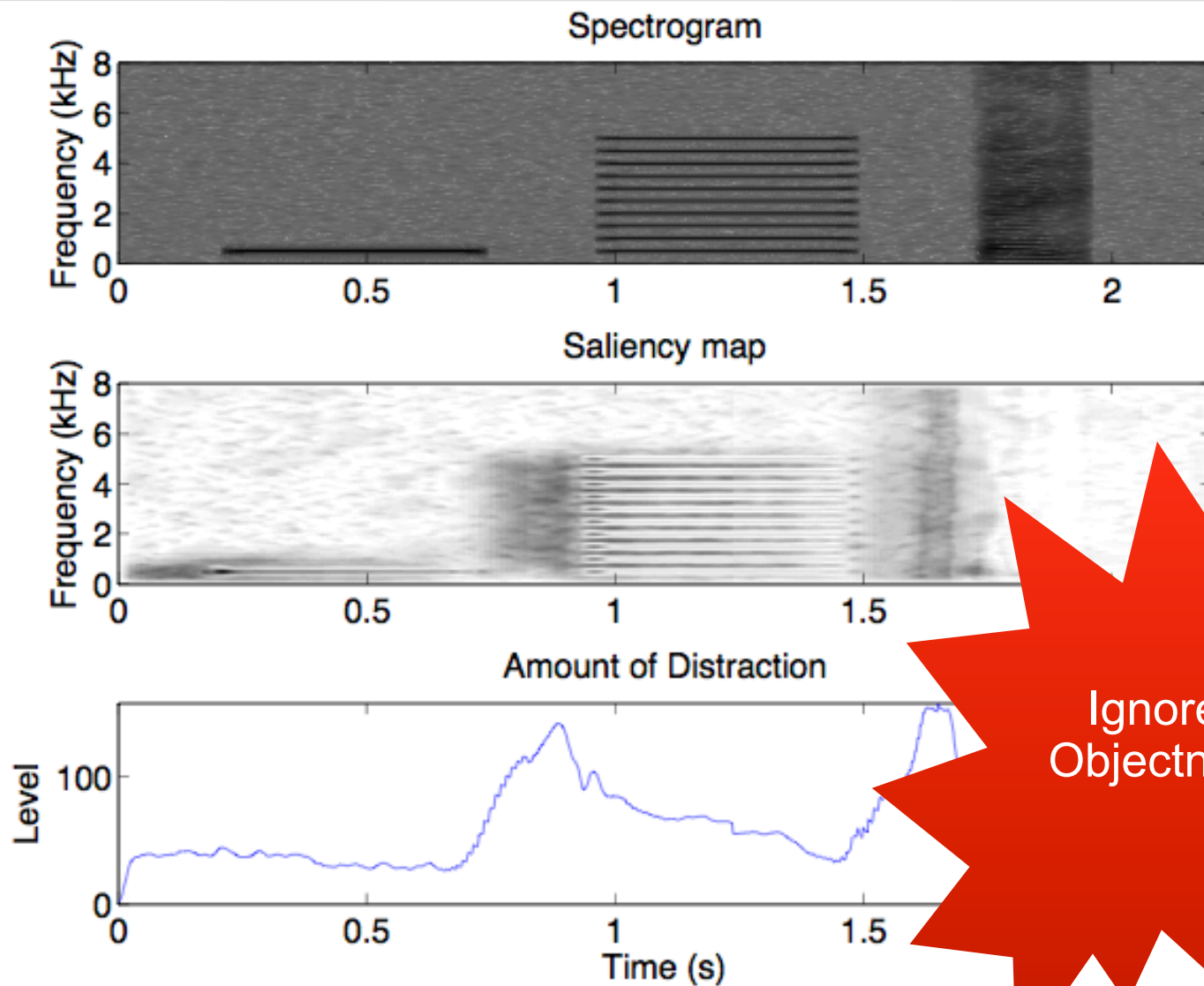
Saliency difference

Contribution of indiv. features

Kayser C, Petkov CI, Lippert M, Logothetis NK. Mechanisms for allocating auditory attention: an auditory saliency map. Curr Biol 2005;15(21):1943-1947.

28

# Kayser Detection Tasks

## Human Tests



**A** Frequency of detection

| [#] | Less Salient | More Salient |
|---|---|---|
| Detected | 206 | 233 |
| Not detected | 82 | 55 |

Fisher's test for interaction of saliency and detection: p<0.01

**B** Detection threshold varies with saliency

Spearman's r: r = 0.56, p<0.01

Detection threshold [dB rel to background]

More noise →

## Monkey Tests



Frequency of detection

| [#] | Less Salient | More Salient |
|---|---|---|
| Detected | 31 | 49 |
| Not detected | 44 | 26 |

Fisher's test for interaction of saliency and detection: p<0.01

Kayser C, Petkov CI, Lippert M, Logothetis NK. Mechanisms for allocating auditory attention: an auditory saliency map. Curr Biol 2005;15(21):1943-1947.

# Kayser Saliency Failures

Spectrogram

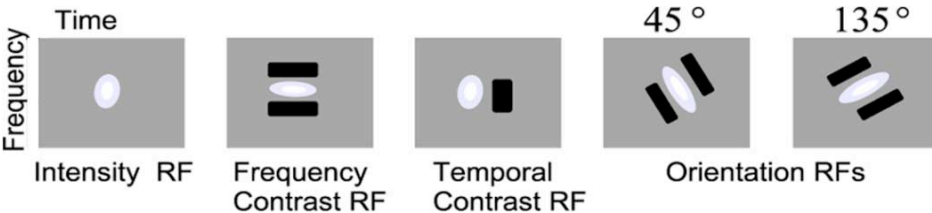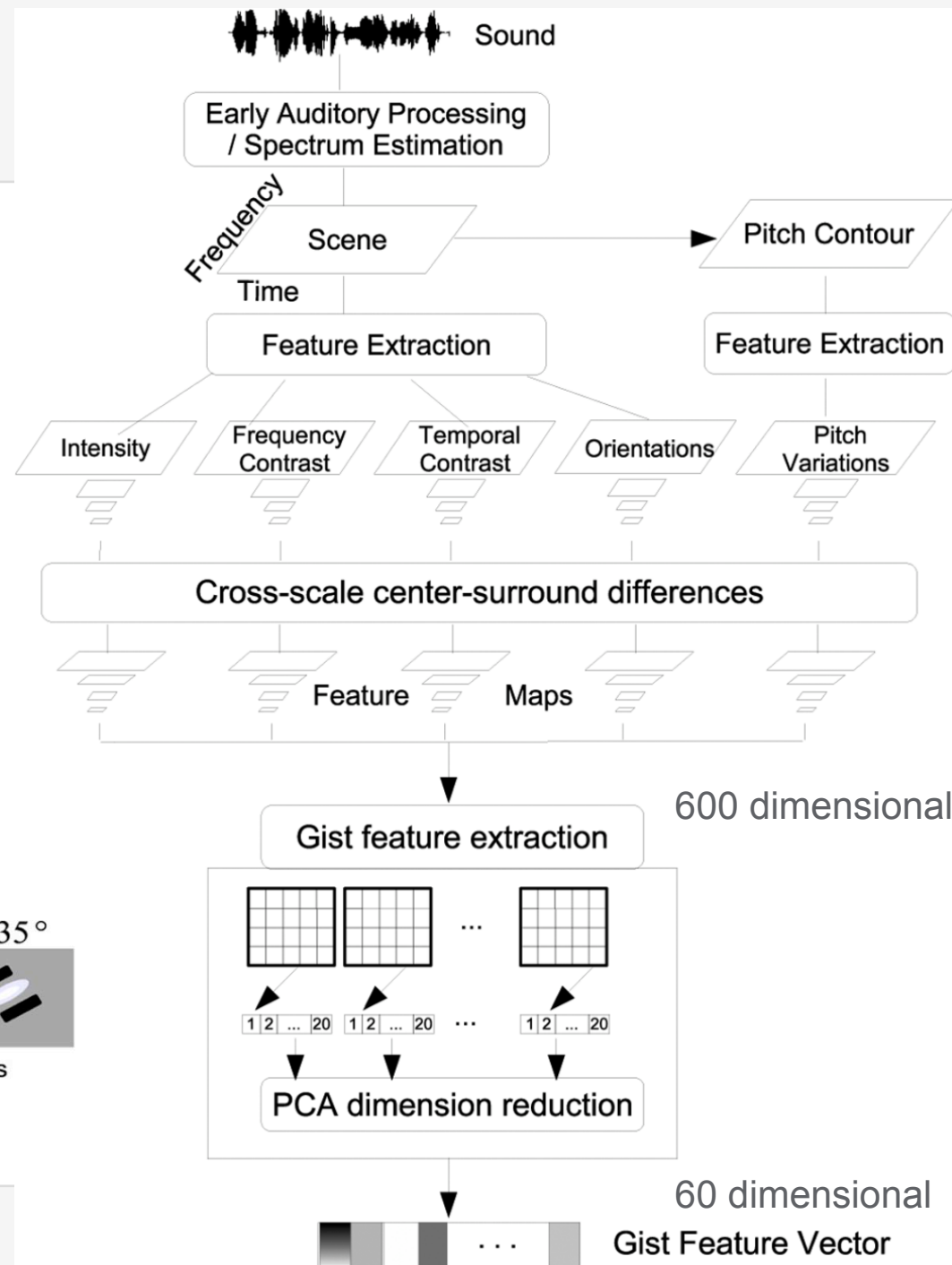Saliency map

Amount of Distraction

Ignores
Objectness

# Kalinli Features

**Extend Itti model**
- Add orientations
- Add pitch variation

**Task**
- Model saliency
- Create gist
- Predict prominence



Kalinili, Ozlem and Narayanan, Shrikanth.A Saliency-Based Auditory Attention Model with Applications to Unsupervised Prominent Syllable Detection in Speech. INTERSPEECH-2007; 2007.

**Gist: "a relatively low-dimensional acoustic scene representation which describes the overall properties of a scene at low-resolution."**



Fig. 3. Auditory attention model. Training phase: the weights are learned in supervised manner. Testing phase: auditory gist features are biased with the learned weights to estimate the top-down model prediction.

Kalinili, Ozlem and Narayanan, Shrikanth.A Saliency-Based Auditory Attention Model with Applications to Unsupervised Prominent Syllable Detection in Speech. INTERSPEECH-2007; 2007.

## Auditory Attention Processing



Utterance Transcription "I'm Irene Doyle".

# Kalinli – Prominence

**Detection**
- All features contribute

**Example analysis**
- Pitch track
- Frequency contrast
- Orientation 45°
- Orientation 135°

**TABLE III**
**PROMINENT SYLLABLE DETECTION PERFORMANCE WITH ONLY PITCH FEATURES**

| Pitch Feature | 1-by-$v$ grids | | | 4-by-5 grids | | |
|---|---|---|---|---|---|---|
| | $d$ | Acc. | F-sc | $d$ | Acc. | F-sc |
| $P_F$ | 21 | 73.90% | 0.57 | 17 | 79.44% | 0.67 |
| $P_{O_{45^o}}$ | 15 | 76.10% | 0.60 | 30 | 79.48% | 0.67 |
| $P_{O_{135^o}}$ | 14 | 74.99% | 0.58 | 29 | 78.65% | 0.65 |
| $P_{O_{45^o}}$ & $P_{O_{135^o}}$ | 26 | 78.88% | 0.66 | 44 | 80.80% | 0.69 |
| $P_F$ & $P_{O_{45^o}}$ & $P_{O_{135^o}}$ | 42 | **80.13%** | 0.68 | 54 | **81.26%** | 0.70 |

(a)

(b)

(c) Detect rising pitch

(d) Detect falling pitch

Fig. 6. Pitch analysis of a speech scene with grid size of 1-by-$v$ (a) pitch. Output obtained with (b) frequency contrast RF. (c) Orientation RF with 45° rotation. (d) Orientation RF with 135° rotation.

# Kalinli – Prominence Detection

## Task

- Detect prominence
- Carry critical information
- Word disambiguation
- More natural TTS

## Features

- Auditory Gist
- Lexical – n-gram prediction
- Syntactic – POS

PROMINENT SYLLABLE DETECTION PERFORMANCE
WITH ONLY PITCH FEATURES

| Pitch Feature | 1-by-$v$ grids | | | 4-by-5 grids | | |
|---|---|---|---|---|---|---|
| | $d$ | Acc. | F-sc | $d$ | Acc. | F-sc |
| $P_F$ | 21 | 73.90% | 0.57 | 17 | 79.44% | 0.67 |
| $P_{O_{45^o}}$ | 15 | 76.10% | 0.60 | 30 | 79.48% | 0.67 |
| $P_{O_{135^o}}$ | 14 | 74.99% | 0.58 | 29 | 78.65% | 0.65 |
| $P_{O_{45^o}}$ & $P_{O_{135^o}}$ | 26 | 78.88% | 0.66 | 44 | 80.80% | 0.69 |
| $P_F$ & $P_{O_{45^o}}$ & $P_{O_{135^o}}$ | 42 | **80.13%** | 0.68 | 54 | **81.26%** | 0.70 |

PROMINENT SYLLABLE DETECTION PERFORMANCE OF INDIVIDUAL
ACOUSTIC, LEXICAL AND SYNTACTIC CUES

| TD Evidence | Acc. | Pr. | Re. | F-sc. |
|---|---|---|---|---|
| Auditory Feat. only | 85.45% | 0.82 | 0.75 | 0.78 |
| Lexical only | 83.85% | 0.77 | 0.76 | 0.76 |
| Syntactic only (word) | 82.50% | 0.82 | 0.87 | 0.84 |
| Syntactic only (syl.) | 68.01% | 0.54 | 0.53 | 0.53 |

COMBINED TOP-DOWN MODEL PERFORMANCE
FOR PROMINENT SYLLABLE DETECTION

| TD Evidence | Acc. | Pr. | Re. | F-sc. |
|---|---|---|---|---|
| Auditory Feat. + Lexical | 88.01% | 0.83 | 0.82 | 0.82 |
| Auditory Feat. + Syntactic | 86.23% | 0.81 | 0.79 | 0.80 |
| Auditory Feat. + Syntactic + Lexical | **88.33%** | 0.83 | 0.83 | 0.83 |
| Combined Feat. word level | **85.71%** | 0.87 | 0.86 | 0.87 |

Kalinili, Ozlem and Narayanan, Shrikanth.A Saliency-Based Auditory Attention Model with Applications to Unsupervised Prominent Syllable Detection in Speech. INTERSPEECH-2007; 2007.

# Duangudom – Modulation Features

## Feature

- Spectrogram
- Spectral-temporal modulation
- Multiscale



Varinthira Duangudom and David Anderson. Using Auditory Saliency To Understand Complex Auditory Scenes. 15th European Signal Processing Conference (EUSIPCO 2007); 2007.

15th European Signal Processing Conference (EUSIPCO 2007), Poznan, Poland, September 3-7, 2007, copyright by EURASIP

# Duangudom – Saliency Experiment

## Model 1

- Scale by $D_i$
- Promotes/inhibits entire feature map

15th European Signal Processing Conference (EUSIPCO 2007), Poznan, Poland, September 3-7, 2007, copyright by EURASIP

Single Peak Enhanced

$M_i$          $M_i^*$

(a)

Multiple peaks suppressed

## Model 2

- Uses local inhibition
- Scales by $D_i$

|  | Correlation to | |
|---|---|---|
| Subject | Model 1 | Model 2 |
| 1 | 0.7324 | 0.8738 |
| 2 | 0.4536 | 0.5329 |
| 3 | 0.4138 | 0.5247 |
| 4 | 0.774 | 0.7872 |
| 5 | 0.0449 | 0.0182 |
| 6 | 0.0632 | 0.1136 |
| 7 | 0.397 | 0.4725 |
| 8 | 0.4073 | |
| 9 | 0.5971 | |
| 10 | 0.5995 | |
| 11 | 0.4234 | 0.4234 |
| 12 | 0.622 | 0.678 |
| 13 | 0.6131 | 0.63 |
| 14 | 0.5447 | 0.5555 |
| Average | 0.4776 | 0.5302 |
| Std Dev | 0.2155 | 0.2379 |

Varinthira Duangudom and David Anderson. Using Auditory Saliency To Understand Complex Auditory Scenes. 15th European Signal Processing Conference (EUSIPCO 2007); 2007.

## Novel features

- Entropy background model

## Inhibition of return

- An orientation mechanism that briefly enhances (for approximately 100–300ms the speed and accuracy with which an object is detected after the object is attended, but then impairs detection speed and accuracy (for approximately 500–3000 ms).



Wang J, Zhang K, Madani K, Sabourin C. Salient environmental sound detection framework for machine awareness. Neurocomputing 2015 3/25;152:444-454.

39

# Wang – Results

**Test sounds**
- Police siren
- Festival with horse steps

**No quantitative comparison**



a
Log power scaled Spectrogram

b
Log power scaled Spectrogram

a
The Auditory Saliency Map Image

b
The Auditory Saliency Map Image

Wang J, Zhang K, Madani K, Sabourin C. Salient environmental sound detection framework for machine awareness. Neurocomputing 2015 3/25;152:444-454.

# Kaya – Temporal Saliency

## Motivation
- Temporal signals
- *Not* images

## Features
- Intensity envelope
- Auditory model
- Lateral inhibition
- STRF
  - Temporal (rate)
  - Bandwidth (spectral ripples)
- Pitch

## Processing
- Multiscale
- Local inhibition
- Threshold
- Sum across channels

Emine Merve Kaya and Mounya Elhilali. A temporal saliency map for modeling auditory attention. Information Sciences and Systems (CISS), 2012 46th Annual Conference on; 2012.

# Kaya – Temporal Saliency Output

## Test signal
- Background: Violin
- Foreground: Flute
- Timing: Frames 450 – 550

## Result with temporal saliency
- Three peaks at:
  Beginning, middle, end

## Comparison to Kayser
- Peaks correspond to background
- No indication near tone



Emine Merve Kaya and Mounya Elhilali. A temporal saliency map for modeling auditory attention. Information Sciences and Systems (CISS), 2012 46th Annual Conference on; 2012.

## Three kinds of test signals

- Timbre: violin->harmonica
- Pitch: 5 semitone rise
- Loudness: 10dB target to mask ratio

## Test

- 20 different variations



| | Our model | Kayser's model | |
|---|---|---|---|
| Hit at 1st peak | 70% | 15% | |
| Hit at 1-3 peaks | 100% | 40% | |
| | **1st peak** | **1st peak** | **1-3 peaks** |
| Hit for timbre | 33.3% | 0% | 0% |
| Hit for pitch | 87.5% | 37.5% | 75% |
| Hit for loudness | 83.3% | 0% | 33% |

Fig. 5.   Detection rates of the target musical notes. Background notes vary only slightly in pitch, while the foreground note can be differing in instrument (timbre), pitch, or loudness. A hit occurs when a peak of the saliency map corresponds to the time of the target note being played.

…nap for modeling auditory attention. Information Sciences and

43

# Kaya – Bottom-up Saliency

## Motivation

- Treat brain as coder
- Predict future (Kalman filter)
- Spike when unexpected
- Focus on intensity, pitch and timbre
- 167D Tensor

Predictive Coding

Envelope

Harmonicity

Spectrogram

Bandwidth

T.Modulation

$x_i(t)$ $x_i'(t)$

$W_{14}$

$W_{25}$

$\alpha(t)$

# Kaya – Bottom-up Saliency Interactions

**Interactions**
- Prior work: Linear sum
- This work: Highly nonlinear



Pitch  Intensity

Harmonicity  Envelope

Spectrogram  Temporal modulation

Bandwidth

Timbre

**B** Contributions of each feature to saliency estimation



Music  Nature  Speech

% of trials

I  H  St  Sb  B  T

# Tsuchida – Auditory Saliency using Natural Statistics

## Saliency definition

$$s_x(t) \propto -\log P(F_x = f_x)$$

- Similar to Bayesian surprise
- Rarity = salience

## Features

- Gammatone
- 20 bands of channels
- PCA
- 2-3 components per band

## Statistics

- GMM with 10 mixtures
- Recent vs. Long past



Audio waveform

Gammatone filterbank

Cochleagram

Split into frequency bands

20

8 ms

PCA   PCA   PCA

Features

Tsuchida T, Cottrell G. Auditory saliency using natural statistics. Society for Neuroscience meeting, New Orleans, LA 2012.

47

**Long tone is more salient**

Spectrogram | Saliency Map



(a) Short and long tones

**Silence is salient in broadband noise**



(b) Gap in broadband noise

**AM modulated tones more salient than stationary**



(c) Stationary and modulated tones

**Second tone less salient**



(d) Single and paired tones

Tsuchida T, Cottrell G. Auditory saliency using natural statistics. Society for Neuroscience meeting, New Orleans, LA 2012.

# Tsuchida – Natural Statistics Results

## Test Material

- Sound effects
- Measure with Kayser and Sun
- 50 high saliency (both)
- 50 low saliency (both)
- 50 large difference (mismatch)
- Subjects

## Tests

- 7 subjects
- 75 pairs

## Compare

Random

Intensity

Kayser

ASUN

# Kim – Machine learned salience

## Training data

- AMI Meeting Corpus
- 12 hours of data
- "Mark the moment when you hear any sound which you unintentionally pay attention to or which attracts your attention."

## Classifier

- Linear on cochlear channel loudness



**Fig. 3.** Estimated impulse response (**h**; Time-Bark plot).

**Table 2**
Equal Error Rate for linear discriminations with the feature combinations.

| Features | Dimension | Equal error rate |
|---|---|---|
| Proposed method (loudness) | 49 | 0.3198 |
| Loudness + zero-crossing rate | 50 | 0.3271 |
| Loudness + spectral flatness | 50 | 0.3958 |
| Loudness + pitch ($T_0$) | 50 | 0.4345 |
| Loudness + $R(T_0)/R(0)$ | 50 | 0.4313 |
| Loudness + all the features | 53 | 0.3922 |
| MFCC | 13 | 0.3446 |

# Saliency – Unsolved Problems

## Data

- No good way to measure saliency effect
- Measuring detectability vs. distraction?
- No common datasets

## Model

- No direct evidence for attention hardware
- Machine-learned vs. Bayesian

# Top-Down

# Attention Changes the Representation?

**Visual Changes with Attention**

Higher

Sharper

**Auditory Changes with Task**



Passive STRF 0:00 · *Discrim* STRF 0:30 · Passive STRF 1:20 · *Discrim* STRF 2:30 · Passive STRF 3:10

Frequency (KHz) — 2, 1, .5, .25

Time (ms) — 0, 40

**Tone detection**

**Chord detection**

**Tone discrimination**

**Voice Activity Detection (VAD)**

**Model adapts**
- Different STRF features for different tasks

**Task?**
- Attention
- Discrimination

Carlin M, Elhilali M. A framework for speech activity detection using adaptive auditory receptive fields. IEEE Trans Acoust , Speech, Signal Process 2015;23(12):2422 - 2433.

# Patil – Task-driven Attentional Mechanism

## Attention modulates
- Sensors
- Object representation

## Task
- Identify one of 12 classes



Patil K, Elhilali M. Task-driven attentional mechanisms for auditory scene recognition. Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP); 2013.

# Mesgarani – Attention Experiment

| Cue | Listen: Two simultaneous speakers | Response |
|---|---|---|
| Ringo | Target : “Ready **Ringo** go to [**Red**] [**Two**] now” | Blue, **Red**, Green / **Two**, Five, Seven |
| | Distractor: “Ready **Tiger** go to [**Green**] [**Five**] now” | |

Target speaker changes randomly from trial to trial.

Target call sign changes after each trial block.

Nima Mesgarani. Personal Communications

Mesgarani N, Fritz J, Shamma S. A computational model of rapid task-related plasticity of auditory cortical receptive fields. J Comput Neurosci 2010 Feb;28(1):19-27.

59

# Mesgarani – Attention Results



$$AMI = Corr\left(SP_1\ attend, SP_1\ alone\right) - Corr\left(SP_1\ attend, SP_2\ alone\right)$$

$$+ Corr\left(SP_2\ attend, SP_2\ alone\right) - Corr\left(SP_2\ attend, SP_1\ alone\right)$$

# Decoding Attention with EEG

# EEG Approaches

- **Single Source (typical BCI)**



Subject is perceiving a dog (or types the letter 'a')

- **Two sources with attention**

20 thousands leagues under the sea

Journey to the center of the earth

Internal Attentional Direction

Subject is attending to "20 thousand"

Attended Signal

Decode and Decide

Where does signal stop?

"Journey to the centre of the earth"

"20,000 leagues under the sea"

# Attention Decoding in a Competing Speaker Environment

# Phrenology?

Speaker 1

Speaker 2

Pick Long-term Winner

State Space Decoder

Average Correlation

Linear Regression

Non-Linear Regression (DNN)

Canonical Correlation

EEG Signals

# Decoding Accuracy

# Time Window

**Correlate**
- Attention decoding accuracy ($r_{attended}$)
- Performance on behavior (memory) task
- r=0.08, P=.005



Subject-Specific | Grand Average

Attended — Unattended — p < 0.05

Low correlation, but significant (blue lines). Tending towards longer reconstruction time lags.

# A Model of Attention

# Scene Analysis Experiments



**Female**

**Male**

Time →

Overlapping Speakers
Two-digit Sentences (even digit at the end)
Template matching (utterance dependent)

Google

**Cognition (Python)**

State: Direction to attend, Digits recognized
Task: Switch attention

HTML (male/female)

HTML (digits)

HTML (Salience)

**ASR (Matlab)**

**Binaural (jAER)**

UDP (ITD)

**Novelty (Python)**

HTML (sound)

# Distracted

- **Switch always**

- **Anytime there is a salient event, switch to active channel**

Male digits are: 98 52 94 34 32 56 14 38 54 94 14 38 58 36 32 38

Female digits are: 54 98 52 12 54 52 56 58 16 14 36 58 16 52 58 52

# Assistive Listening

- **Speech vs. music**

- **Different speakers**

- **Reduce environmental noise**



Selective Hearing For Men

# Speech Recognition

## With eye gaze

"For all the talk among Democrats"



Time and radius parameters?

# Modifying Language Models

- **Ideal system**
- **Current implementation**

Generic + gaze-specific language model →

[ ASR ] →

Generic + page-specific language model →

[ ASR ] — N-best lattice →

Gaze-specific language model →

[ Rescoring ] →

# Lattice Rescoring

- **Get recognition results**
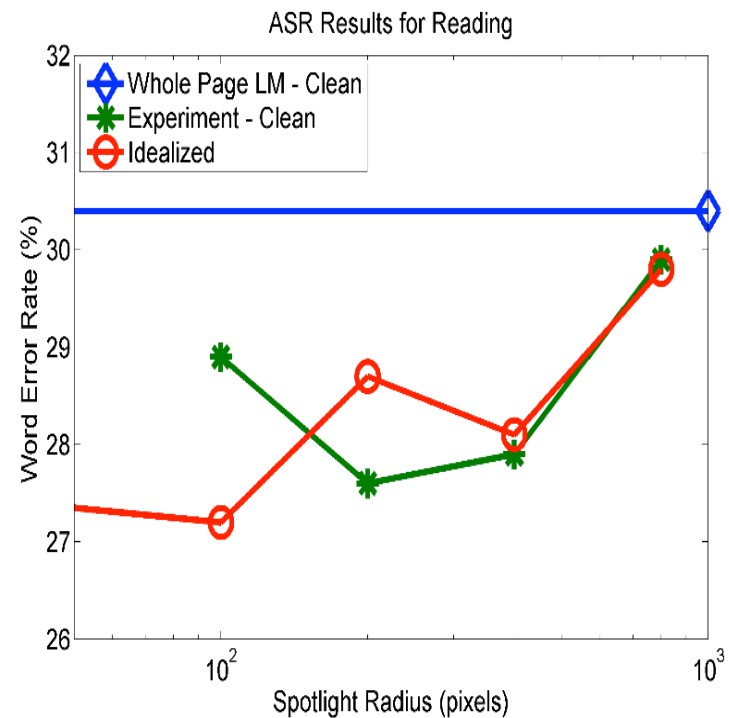  Rescore transition probabilities
  First pass: This is a test sentence.
  Second pass (with eye gaze): This is the guest sense.
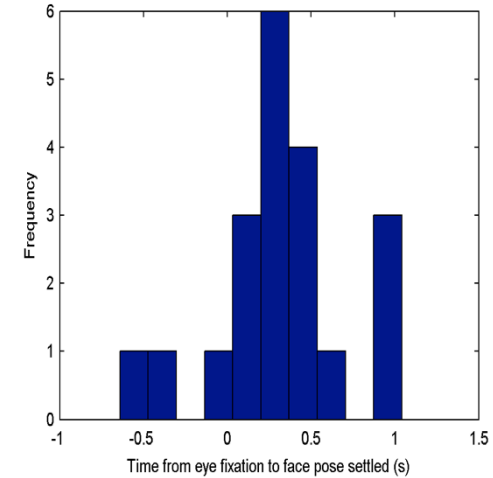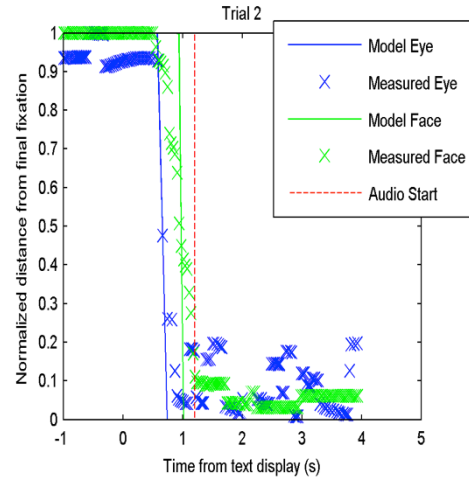
# ASR with Eye Gaze

- **Using eye gaze reduces LM perplexity**

| Language Model | Perplexity |
|---|---|
| Generic (GLM) | >1000 |
| GLM + page | 26 |
| GLM + gaze | 15 |
| GLM + page + gaze | 14 |

- **Approximately a 10% relative error rate reduction**
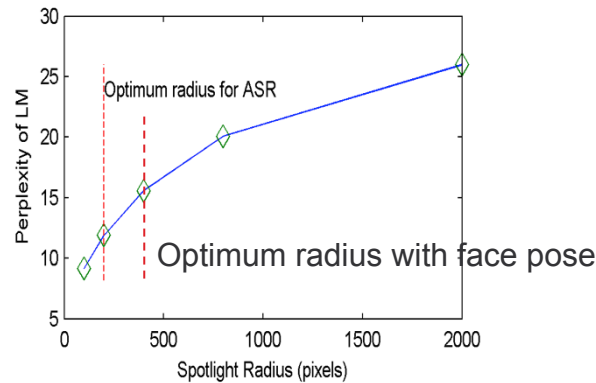


ASR Results for Reading

# Face Pose Approximates Eye Gaze

- **Timing**

- **Radius**



Radius doubles, but perplexity only goes up by ~30%

# Demo

# Gaze-Enhanced Speech Recognition

**Please say one of the phrases from the boxes on the right.**

## Context #1

| Sent email | Do no evil | Microsoft Office | Wreck a nice beach |
|---|---|---|---|

## Context #2

| Send email | Do you know evil? | Microsoft Office | Recognize speech |
|---|---|---|---|

# Conclusions

**Salience matters**

**What's the right model?**

**Where do we get data?**

**What kind of data?**

# Thank you

malcolm@ieee.org

**Task**

**Measure STRF of neurons**

**Change before and after**

Yin P, Fritz JB, Shamma SA. Rapid spectrotemporal plasticity in primary auditory cortex during behavior. J Neurosci 2014 Mar 19;34(12):4396-4408.

86