# Enhancement of Noisy Speech

# State-of-the-Art and Perspectives

Rainer Martin

Institute of Communications Technology (IFN)

Technical University of Braunschweig

July, 2003

# Applications of Noise Reduction

- Hands-free telephony.

- Robust speech recognition.

- Robust speech coding (ETSI/3GPP AMR, MELPe, ITU-T 4 kbit/s codecs).

- Hearing aids and cochlear implants.

- Restoration of historic recordings.

- Forensic applications.

# Ingredients

- Models of speech production

- Signal theory

- Room acoustics

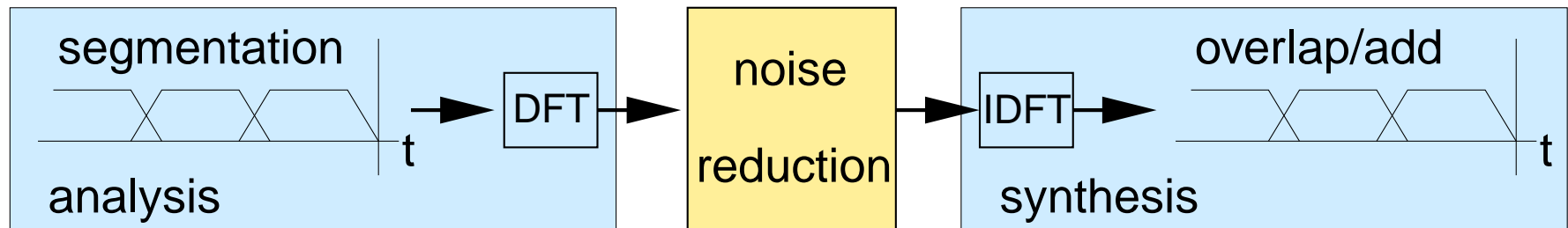- Psychoacoustics

- Models of speech perception

<p style="text-align:center; color:blue">Objective: Improve quality <u>and</u> intelligibility!</p>

<p style="text-align:center; color:red">Combine signal theoretic and perceptive approaches!</p>

# Noise Reduction in the Spectral Domain

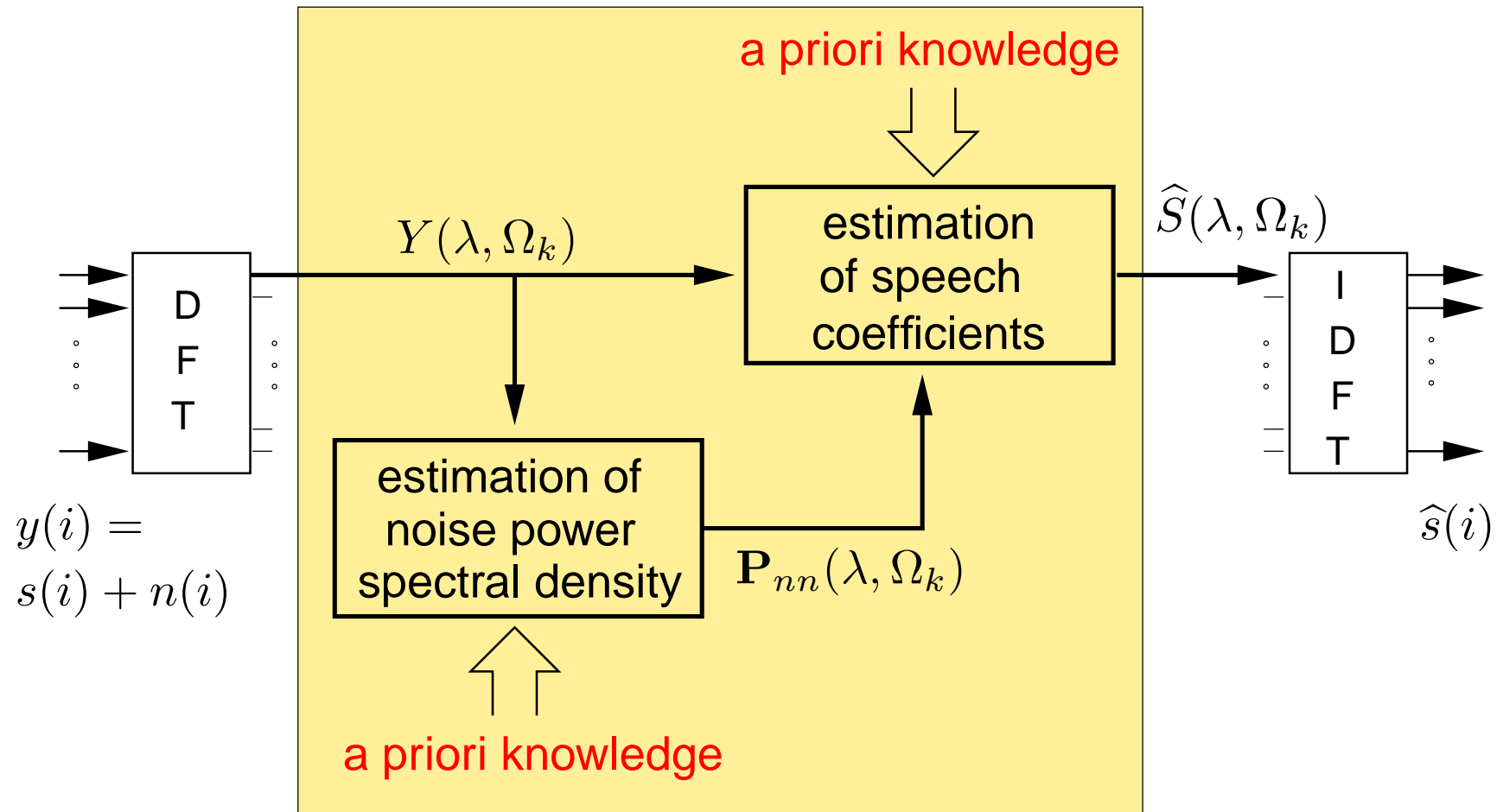▶ **Spectral analysis – noise reduction – synthesis:**



▶ **Advantages of spectral processing:**

- good separation of speech and noise

- decorrelation of spectral components

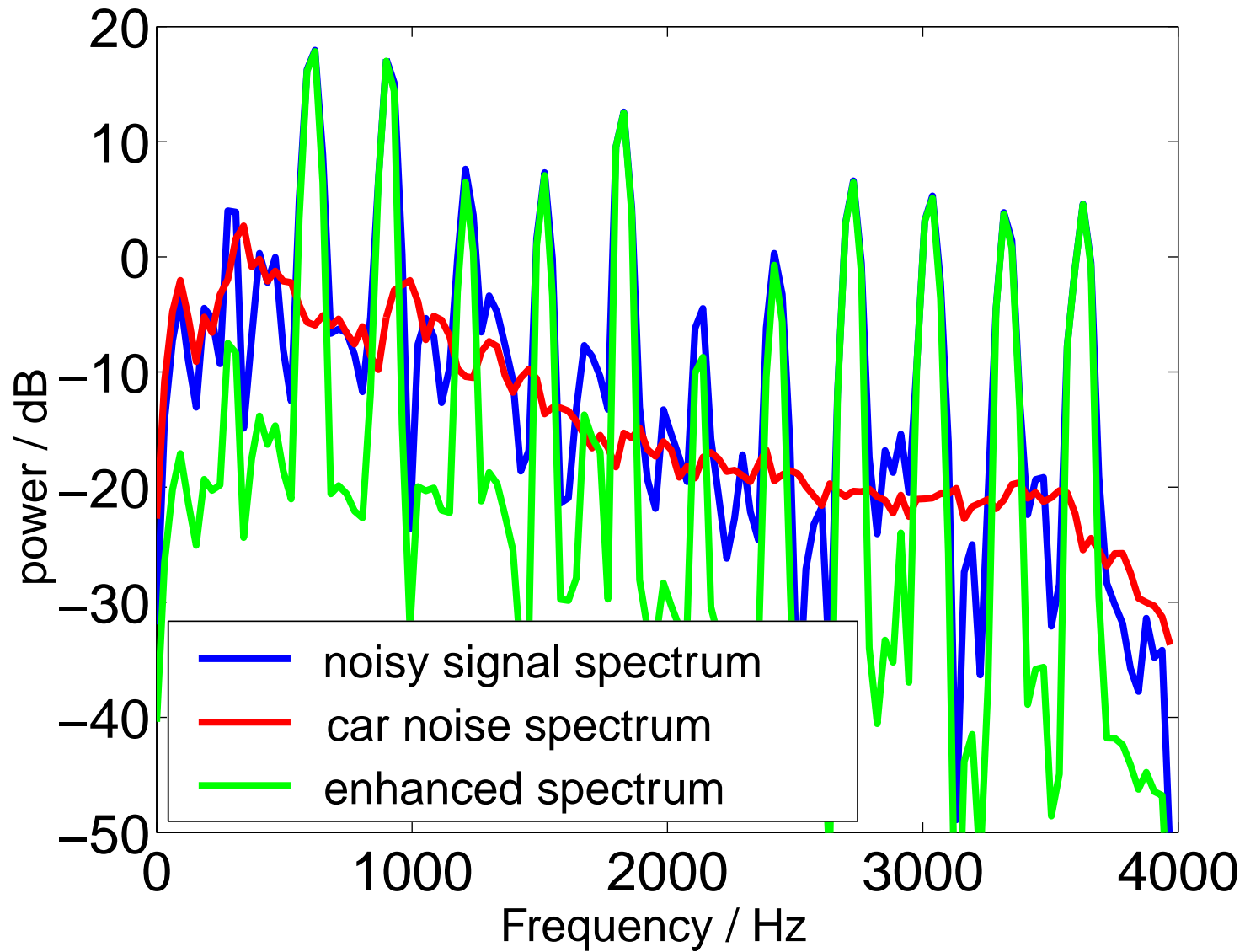- integration of psychoacoustic models

# Principles of Noise Reduction

$\lambda \longrightarrow$ frame index $\qquad k \longrightarrow$ frequency bin index

# Principles of Noise Reduction



placements

ise spectrum

ced spectrum

# Estimation of Speech Coefficients

▶ **Linear estimators**

- e.g. Wiener Filter

▶ **Non-linear estimators**

- MMSE Short Time Spectral Amplitude estimator

  [Ephraim & Malah, 1984, 1985]

- Psychoacoustic methods [Gustafsson et al. 1998]

- MMSE estimation based on supergaussian priors

  [Martin 2002]

# MMSE Estimation

► **Optimal estimate for independent real and imaginary parts:**

$$E\{S \mid Y\} = E\{S_R \mid Y_R\} + jE\{S_I \mid Y_I\}$$

► **Estimation of either the real or the imaginary part:**

$$E\{S_\diamond \mid Y_\diamond\} = \int_{-\infty}^{\infty} S_\diamond p(S_\diamond \mid Y_\diamond)dS_\diamond$$

► **Application of Bayes theorem:**

$$E\{S_\diamond \mid Y_\diamond\} = \frac{1}{p(Y_\diamond)} \int_{-\infty}^{\infty} S_\diamond p(Y_\diamond \mid S_\diamond)p(S_\diamond)dS_\diamond$$

► **What is the appropriate prior density $p(S_\diamond)$ ?**

# Some Answers and Some Questions

▶ **DFT coefficients are asymptotically complex Gaussian distributed ! [Brillinger, 1981]**

▶ **Typical frame size in mobile communications: 10-30 ms $<$ span of correlation of (voiced) speech !**

▶ **Do the asymptotic assumptions hold for speech signals ???**

▶ **No! See, e.g., [Porter and Boll, 1984].**

# Prior Densities for Real and Imaginary Part

▶ **Gaussian pdf:**

$$p(S_\diamond) = \frac{1}{\sqrt{\pi}\sigma_s} \exp\left(-\frac{S_\diamond^2}{\sigma_s^2}\right)$$

$\rightarrow$ Wiener filter

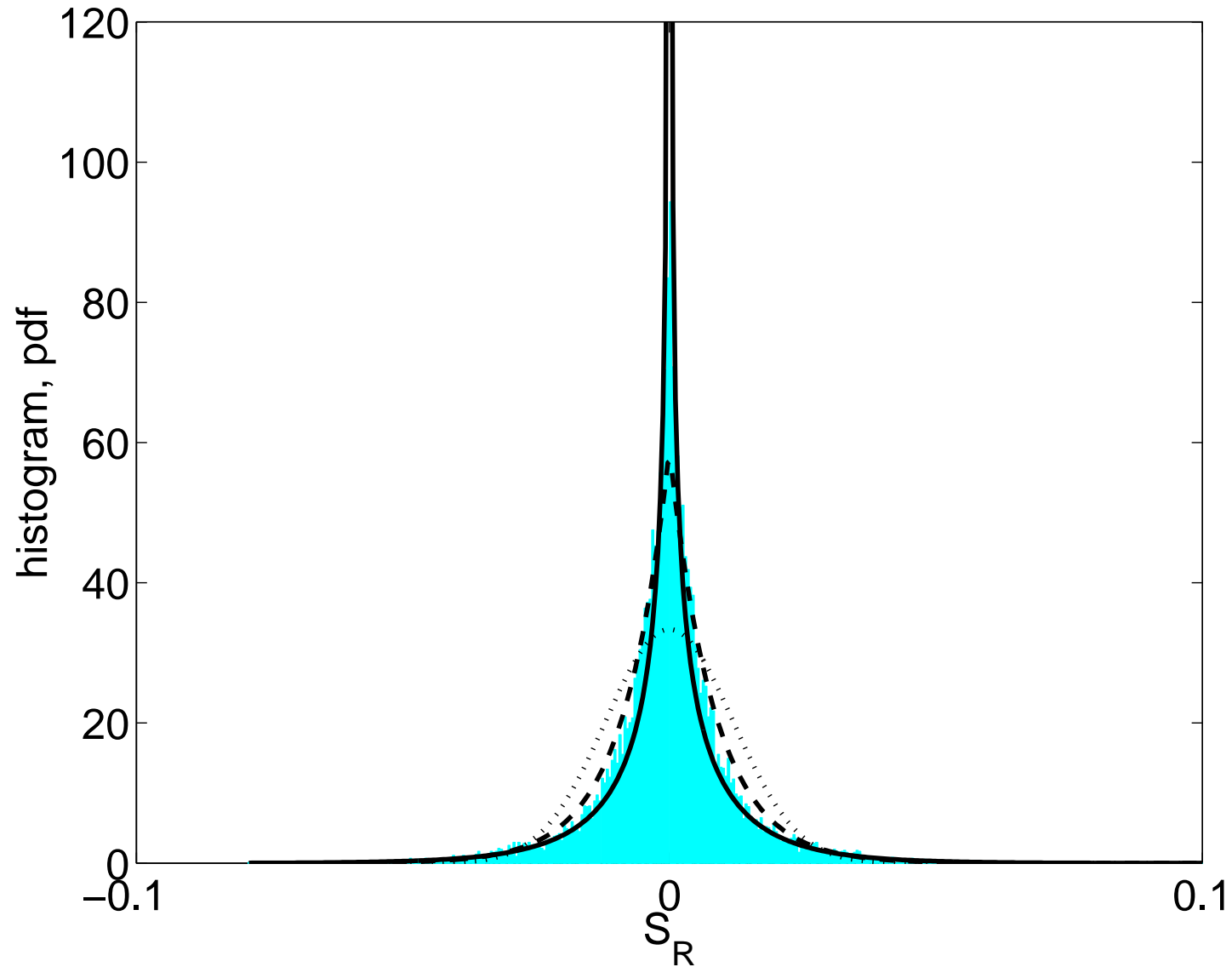▶ **Laplacian pdf:**

$$p(S_\diamond) = \frac{1}{\sigma_s} \exp\left(-\frac{2|S_\diamond|}{\sigma_s}\right)$$

▶ **Gamma pdf:**

$$p(S_\diamond) = \frac{\sqrt[4]{3}}{2\sqrt{\pi\sigma_s}\sqrt[4]{2}}|S_\diamond|^{-\frac{1}{2}} \exp\left(-\frac{\sqrt{3}|S_\diamond|}{\sqrt{2}\sigma_s}\right)$$

# Histogram of DFT Coefficients for Speech
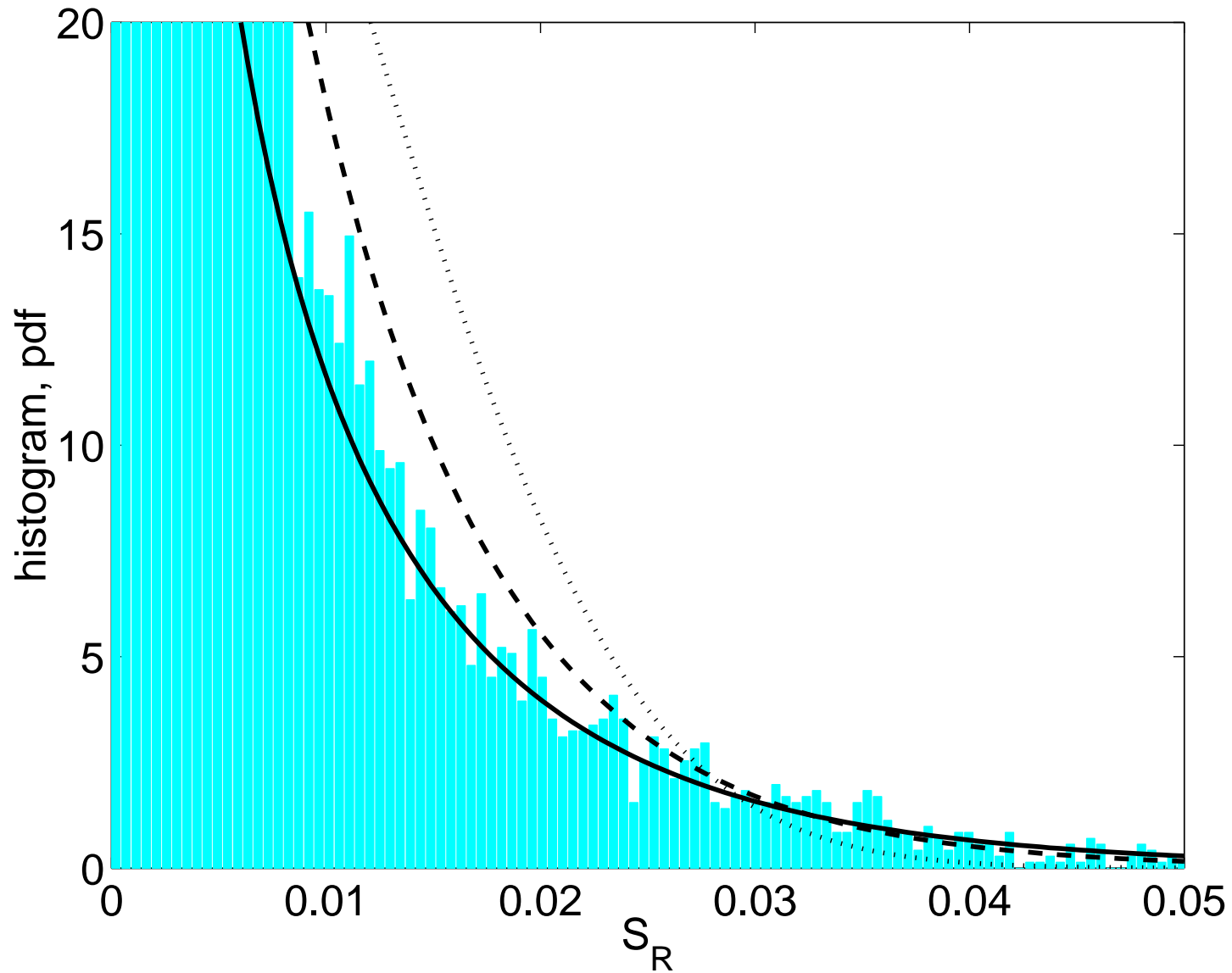


dotted: Gaussian pdf          dashed: Laplacian pdf          solid: Gamma pdf

# Histogram of Speech Coefficients (enlarged)



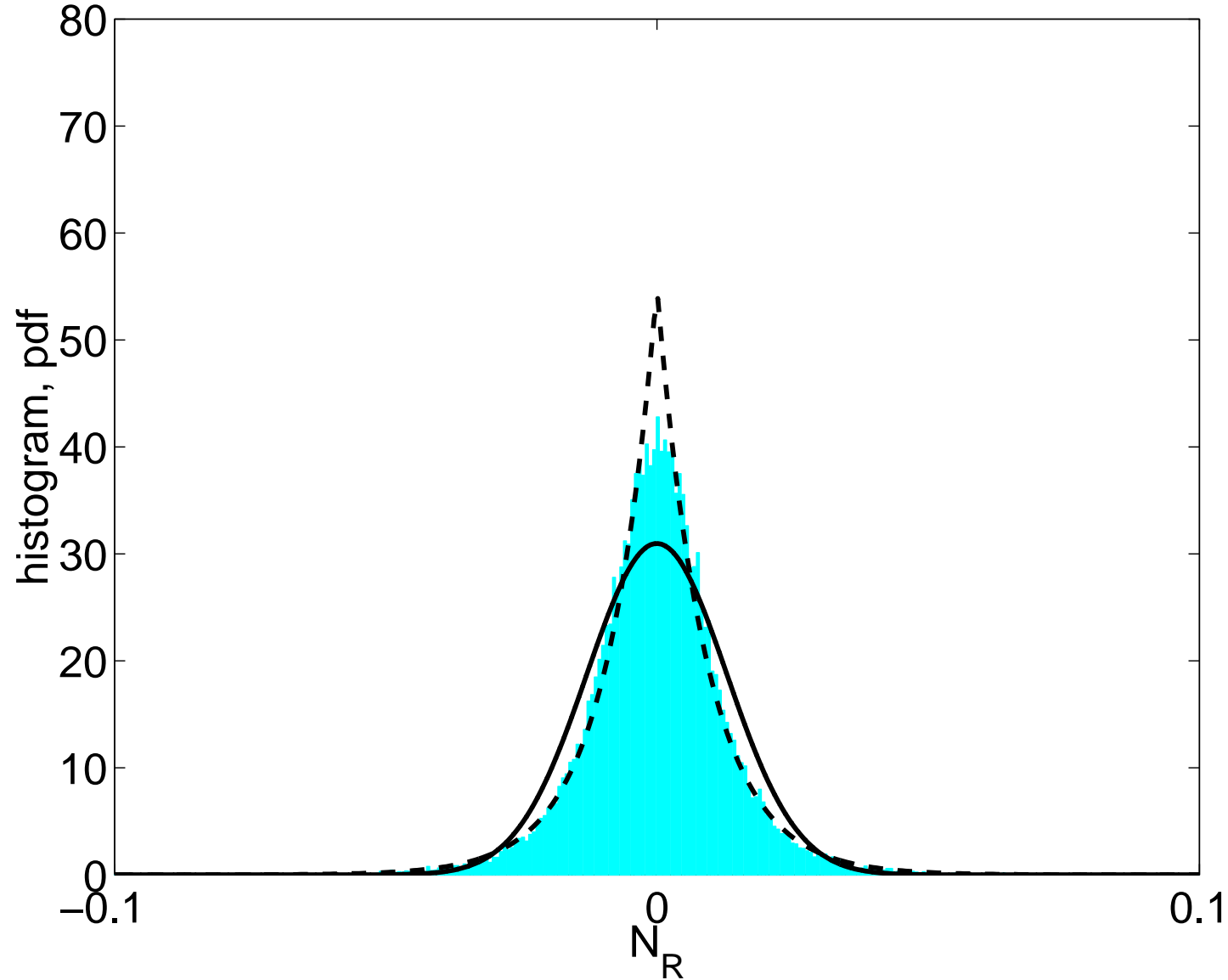dotted: Gaussian pdf      dashed: Laplacian pdf      solid: Gamma pdf

# Histogram of DFT Coefficients for Car Noise



dotted: Gaussian pdf        dashed: Laplacian pdf

# Histogram of Car Coefficients (enlarged)



dotted: Gaussian pdf          dashed: Laplacian pdf

# Non-linear MMSE Estimator
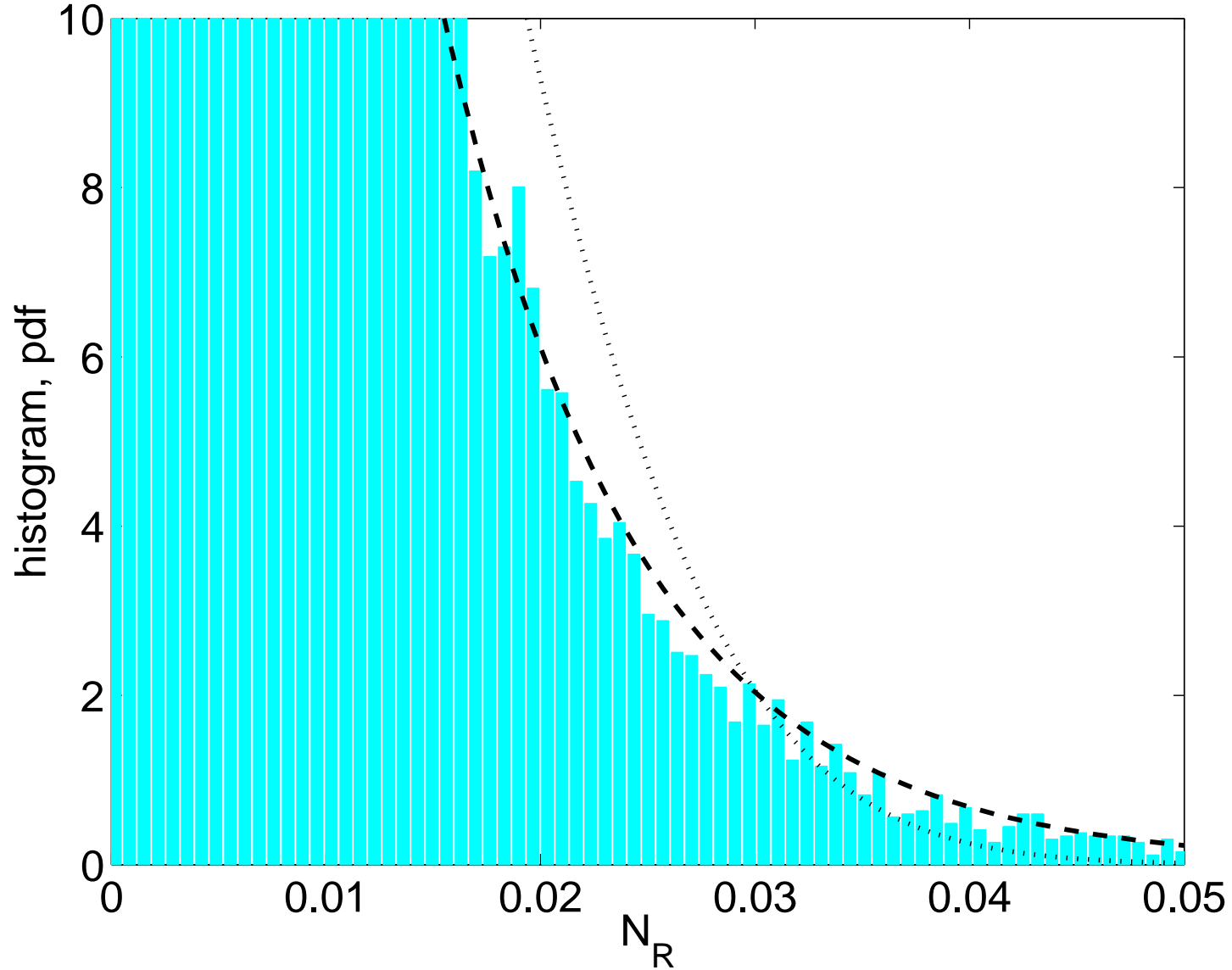


Laplacian Noise and Gamma Speech Prior

$$\sigma_s^2 + \sigma_n^2 = 2$$

# Segmental SNR Improvement (White Noise)

# Relative Improvement w.r.t. Wiener Filter



ag replacements

segmental SNR of input signal

# Background Noise PSD Estimation

▶ **Methods:**

- Voice activity detection;

- Soft-decision methods;

- Biased compensated tracking of spectral minima

  [Martin 1994, 2001]

▶ **Assumptions:**

- Speech and noise are statistically independent;

- Speech is not always present;

- Noise is more stationary than speech.

# Minimum Statistics: Basic Principle

# Minimum Statistics: Bias

# Mean of Minimum

$B_{min}^{-1}$



150

E{minimum}

$Q_{eq} = 512$

$Q_{eq} = 128$

$Q_{eq} = 64$

$Q_{eq} = 32$

$Q_{eq} = 16$

$Q_{eq} = 8$

$Q_{eq} = 4$

$Q_{eq} = 2$

$q = 256$

$D$

$D$: length of minimum search window

$Q = 64$

$Q_{eq} = 1/var\{P(\lambda, \Omega_k)\}_{norm}$

$Q = 128$

$Q = 256$

$Q = 512$

fN

# Minimum Statistics: What's New ?

$B_{min}^{-1}$

$Q = 2$

$Q = 4$

$Q = 8$

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$

▶ **Minimum Statistic, version 1994**

- fixed smoothing parameter $\alpha$

- fixed bias compensation

▶ **Minimum Statistic, version 2001**

- signal dependent optimal smoothing $Q = 2$

- signal dependent bias compensation $Q = 8$

$Q = 4$

- fast minimum update

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$

ffN

# Minimum Statistics (version 2001)

$Q = 128$
$Q = 256$ $Q = 128$
$Q = 512$ $Q = 256$
$Q = 2$ $Q = 512$
$Q = 4$ $Q = 2$
$Q = 8$ $Q = 4$
$Q = 16$ $Q = 8$
$Q = 32$ $Q = 16$
$Q = 64$ $Q = 32$
$Q = 128$ $Q = 64$
$Q = 256$ $Q = 128$
$Q = 512$ $Q = 256$
$Q = 512$

- periodogram (frequency bin k=25)
- smoothed periodogram (k=25)
- minimum of smoothed periodogram

90

80

70

dB 60

50

40

dB 30

frequency bin k=25)
periodogram (k=25)
smoothed periodogram

frame index

Estimation of noise power spectral density without voice activity detection !

$Q = 256$
$Q = 512$

# Relative Estimation Error

$B_{min}^{-1}$

$Q = 2$

$Q = 4$

▶ **Speech pause:**

$Q = 8$

| Algorithms | white noise | vehicular noise | street noise |
|---|---|---|---|
| MinStat 1994 ($\alpha = 0.6$) | 0.059 (0.11) | 0.062 (0.13) | -0.15 (0.21) |
| MinStat 2001 | -0.006 (0.041) | -0.016 (0.041) | -0.27 (0.13) |

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

(in parentheses: variance of estimation error)

$Q = 256$

$Q = 512$

▶ **Speech activity (3 min without speech pauses):**

$Q = 2$

$Q = 4$

| Algorithms | white noise | vehicular noise | street noise |
|---|---|---|---|
| MinStat 1994 ($\alpha = 0.6$) | 0.64 (0.77) | 0.77 (1.04) | 0.59 (1.9) |
| MinStat 2001 | -0.04 (0.14) | 0.02 (0.17) | -0.20 (0.28) |

$Q = 8$

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$

# Two Channel Noise Reduction

$B_{min}^{-1}$

$Q = 2$

$Q = 4$

$Q = 8$

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$

$x_1(k)$

preem-phasis

$T$

$\vec{h}_1$

$T_H$

adaptive time delay estimation

$T_H$

$\vec{h}_2$

$+$  $-$

$+$  $+$  $-$

$\frac{\vec{h}_1 + \vec{h}_2}{4} \otimes \vec{w}$

deem-phasis

$y_{\mathrm{hppre1}}(k)$

$y_{\mathrm{hppre2}}(k)$

$x_2(k)$

preem-phasis

$\Delta T$

$Q = 2$

$Q = 4$

$Q = 8$

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$
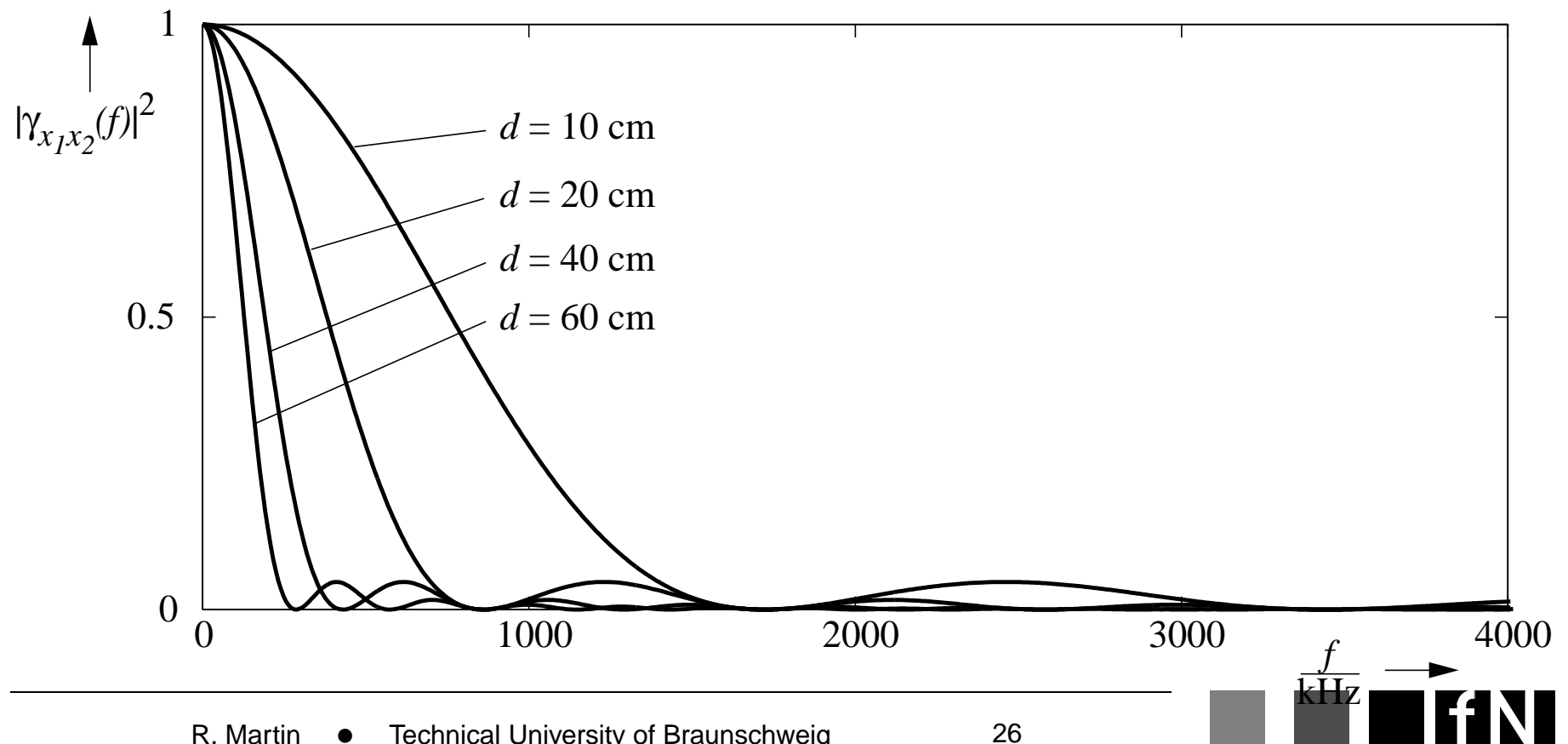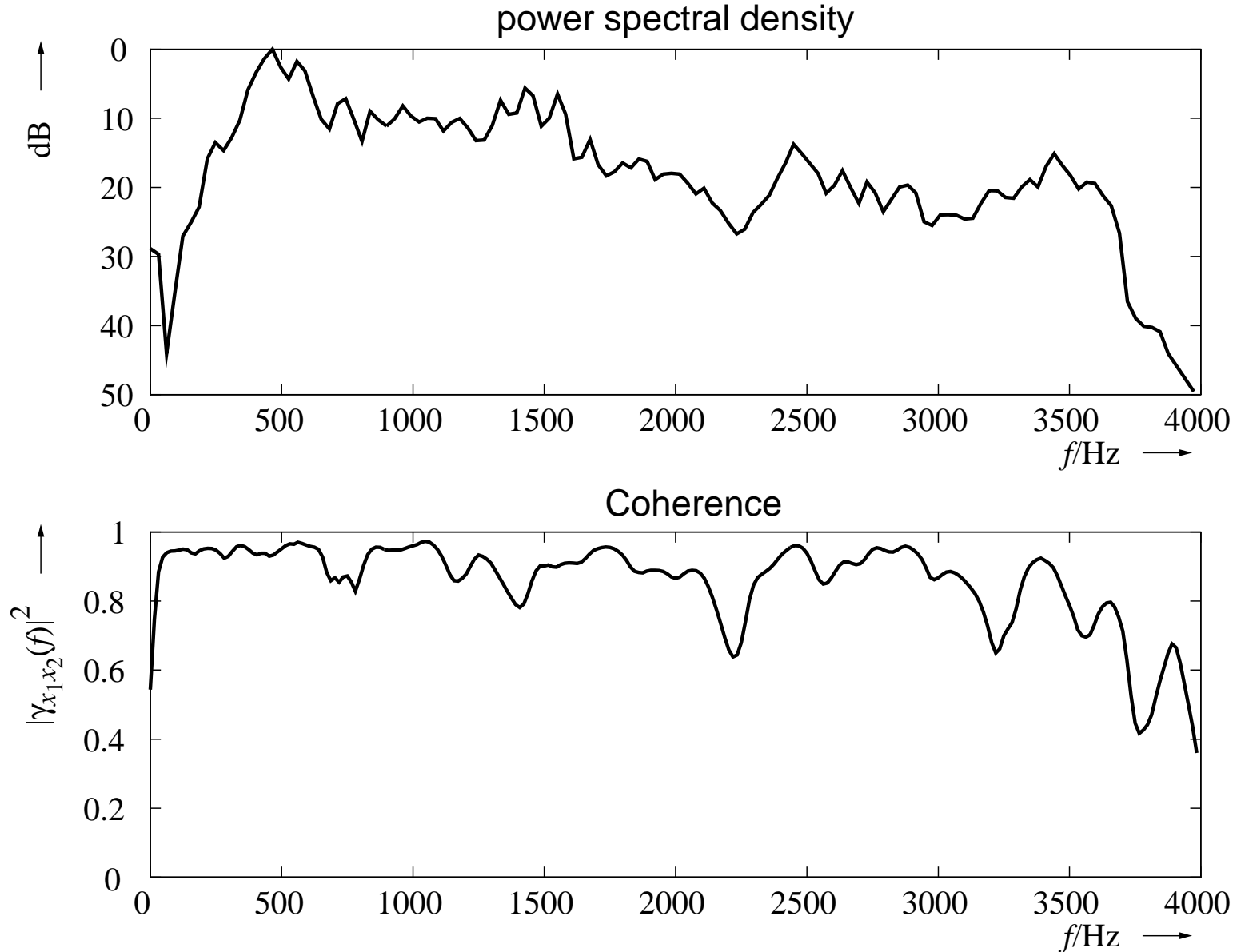
# Coherence of Noise (Diffuse Sound Field)

The complex coherence $\gamma_{x_1x_2}(\Omega)$ of two signals $x_1(k)$ and $x_2(k)$ is defined as

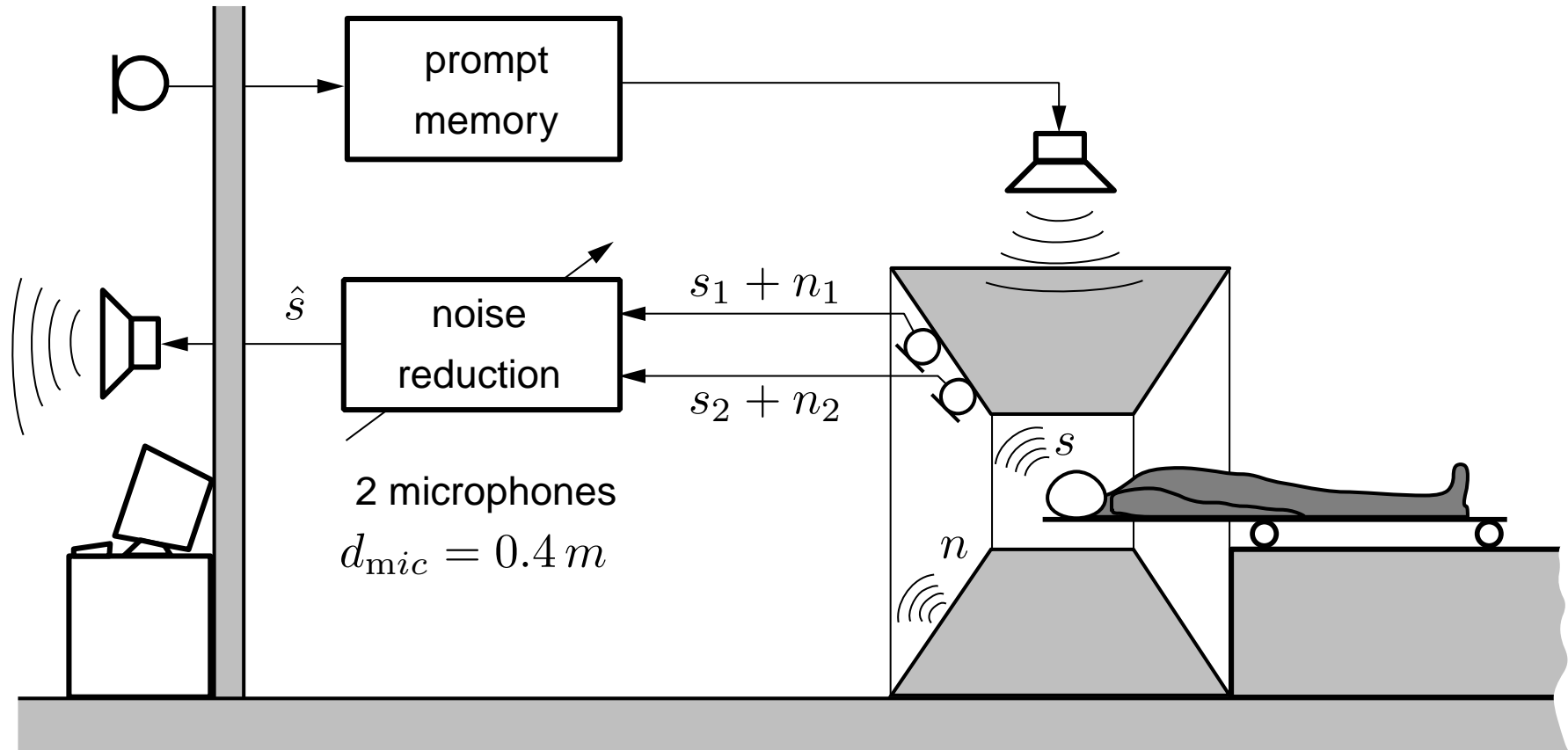$$\gamma_{x_1x_2}(\Omega) = \frac{\Phi_{x_1x_2}(e^{j\Omega})}{\sqrt{\Phi_{x_1x_1}(e^{j\Omega})\,\Phi_{x_2x_2}(e^{j\Omega})}} \ .$$

$d = 10$ cm
$d = 20$ cm
$d = 40$ cm
$d = 60$ cm

$|\gamma_{x_1x_2}(f)|^2$

$\frac{f}{\text{kHz}}$

# Coherence of Speech in a Car

$Q = 2$

$Q = 4$

$Q = 8$



power spectral density

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

dB

$f$/Hz

$Q = 2$

$Q = 4$

$Q = 8$

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$

$Q = 512$

Coherence

$|\gamma_{x_1 x_2}(f)|^2$

$Q = 2$

$Q = 4$

$Q = 8$

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$f$/Hz

$Q = 256$

$Q = 512$

fN

# Two Channel Noise Reduction

$B_{min}^{-1}$

$Q = 2$

$Q = 4$

$Q = 8$

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$

prompt
memory

noise
reduction

$\hat{s}$

$s_1 + n_1$

$s_2 + n_2$

2 microphones

$d_{\mathrm{mic}} = 0.4\,m$

$s$

$n$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$

# First-Order Differential Microfone

$Q = 2$
$Q = 4$
$Q = 8$
$Q = 16$
$Q = 32$
$Q = 64$
$Q = 128$
$Q = 256$
$Q = 512$

$B_{min}^{-1}$

$Q = 2$
$Q = 4$
$Q = 8$
$Q = 16$
$Q = 32$
$Q = 64$
$Q = 128$
$Q = 256$
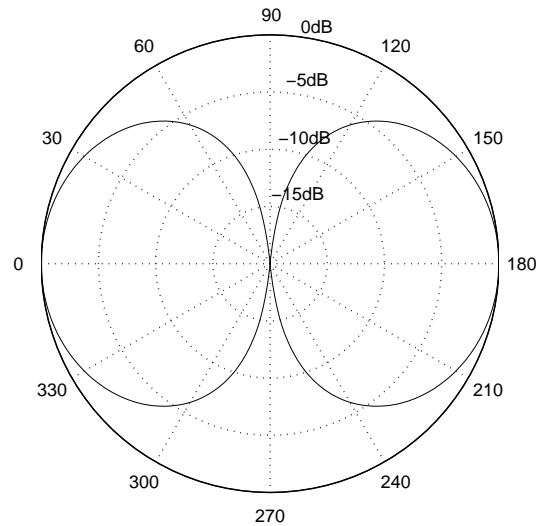$Q = 512$

Speech

d

$\alpha$

Delay T

$-$ $+$

Equalization

$$Y(j\omega) = S(j\omega)e^{j\omega\left(\frac{d}{2c}\cos(\alpha)\right)}\left[1 - e^{-j\omega\frac{d}{c}\left(\cos(\alpha)+\frac{cT}{d}\right)}\right]$$

$$\left|\frac{Y(j\omega)}{S(j\omega)}\right| = 2\left|\sin\left(\frac{\omega d}{2c}\left(\cos(\alpha)+\frac{cT}{d}\right)\right)\right|$$
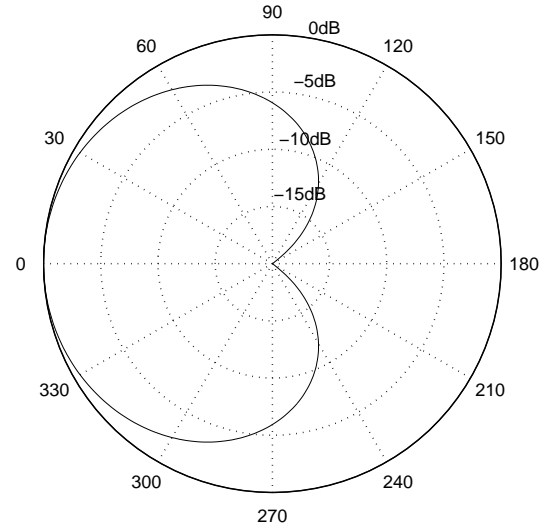
$Q = 2$
$Q = 8$
$Q = 16$
$Q = 32$
$Q = 64$
$Q = 128$
$Q = 256$
$Q = 512$

# Directivity Patterns ($d = 0.015$ m, $f = 1kHz$ )

Dipole ( Tc/d = 0), f = 1000 Hz



Cardioid ( Tc/d = 1), f = 1000 Hz

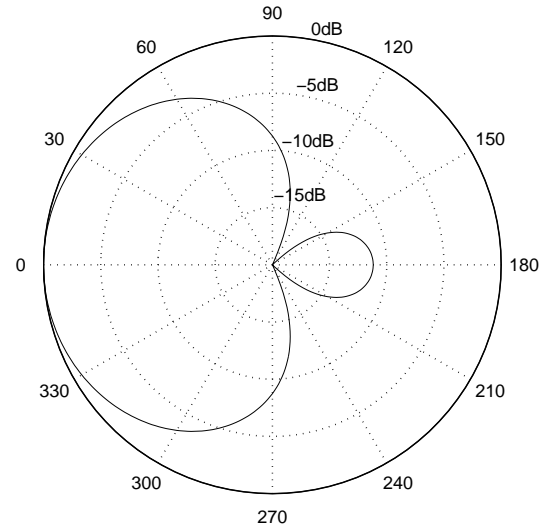Azimuth angle in degrees

Hyper Cardioid ( Tc/d = 0.34), f = 1000 Hz

Azimuth angle in degrees



Super Cardioid ( Tc/d = 0.57), f = 1000 Hz

Azimuth angle in degrees

$Q = 256$

$Q = 512$

# Delay-and-Sum Beamformer

$B_{min}^{-1}$

$Q = 2$



source

$s(k)$

$\theta$

$y_1(k)$ $T_1$ $\widetilde{y}_1(k)$

$y_2(k)$ $T_2$ $\widetilde{y}_2(k)$

$y_3(k)$ $T_3$ $\widetilde{y}_3(k)$

$y_N(k)$ $T_N$ $\widetilde{y}_N(k)$

$\widehat{y}(k)$

noise $n_\ell(k)$

i.i.d. noise: Gain $G = 10\log(N)$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$

PSfrag replacements

$Q = 2$

$Q = 4$

$B_{min}^{-1}$

$Q = 8$

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$

$Q = 2$

$Q = 4$

$Q = 8$

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$



$Q = 256$

$Q = 512$

# Directivity Pattern

PSfrag replacements

$Q = 2$
$Q = 4$
$Q = 8$
$Q = 16$
$Q = 32$
$Q = 64$
$Q = 128$
$Q = 256$
$Q = 512$
$Q = 2$
$Q = 4$
$Q = 8$
$Q = 16$
$Q = 32$
$Q = 64$
$Q = 128$
$Q = 256$
$Q = 512$

$B_{min}^{-1}$

$Q = 128$

$Q = 256$

$Q = 512$

# Arrays for Speech Acquisition in Cars

$Q = 256$
$Q = 512$

$B-1$
$B_{min}$

$Q = 2$            $Q = 2$

$Q = 4$

$Q = 8$

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$



● 1    ● 2    ● 3    ● 4    5 ●

5 cm    4 cm    4 cm    5 cm

5.25 cm    X    Y

● 8    ● 7    6 ●

$Q = 2$

- microphones 1, 2, 3, 4, 5 → linear array $Q = 4$

$Q = 8$

- microphones 1, 2, 7, 4, 5 → planar array $Q = 16$

$Q = 32$

$Q = 64$ [Martin et al. 2001]

$Q = 128$

$Q = 256$

# Delay-and-sum vs. Superdirective Arrays

# Linear and Planar Microphone Arrays

# Adaptive Beamformer (GSC)

# Conclusions

$B_{min}^{-1}$

$Q = 2$

$Q = 4$

$Q = 8$

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$

▶ **Find better ways to exploit statistics of signals!**

● Incorporate models of speech production

● Develop better background noise estimation methods

● Design algorithms for high quality and intelligibility

● Exploit spatial selectivity using multiple microphones

$Q = 2$

$Q = 4$

$Q = 8$

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$

▶ **Understand processing in the auditory system:**

● Enhance perceptionally important features

● Use perceptive models to reduce complexity of algorithms

# Selected References

$$B_{min}^{-1}$$

$Q = 2$

$Q = 4$

$Q = 8$

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$

$Q = 2$

$Q = 4$

$Q = 8$

$Q = 16$

$Q = 32$

$Q = 64$

$Q = 128$

$Q = 256$

$Q = 512$