

Perception and Binaural Signal Processing

EE4715 Array Processing | 7 June '24

Audio Signal Processing Applications



MVDR Beamforming

- STFT domain (short-time stationarity assumption)
- narrowband assumption
- minimise output noise power
- target preservation constraint

$$\mathbf{x}(k, l) = \mathbf{a}(k, l)s(k, l) + \mathbf{n}(k, l)$$

$$y(k, l) = \mathbf{w}^H(k, l)\mathbf{x}(k, l)$$

$$\text{minimise } \mathbf{w}^H(k, l)\mathbf{R}_n(k, l)\mathbf{w}(k, l)$$
$$\mathbf{w}(k, l)$$

$$\text{subject to } \mathbf{w}^H(k, l)\mathbf{a}(k, l) = 1$$

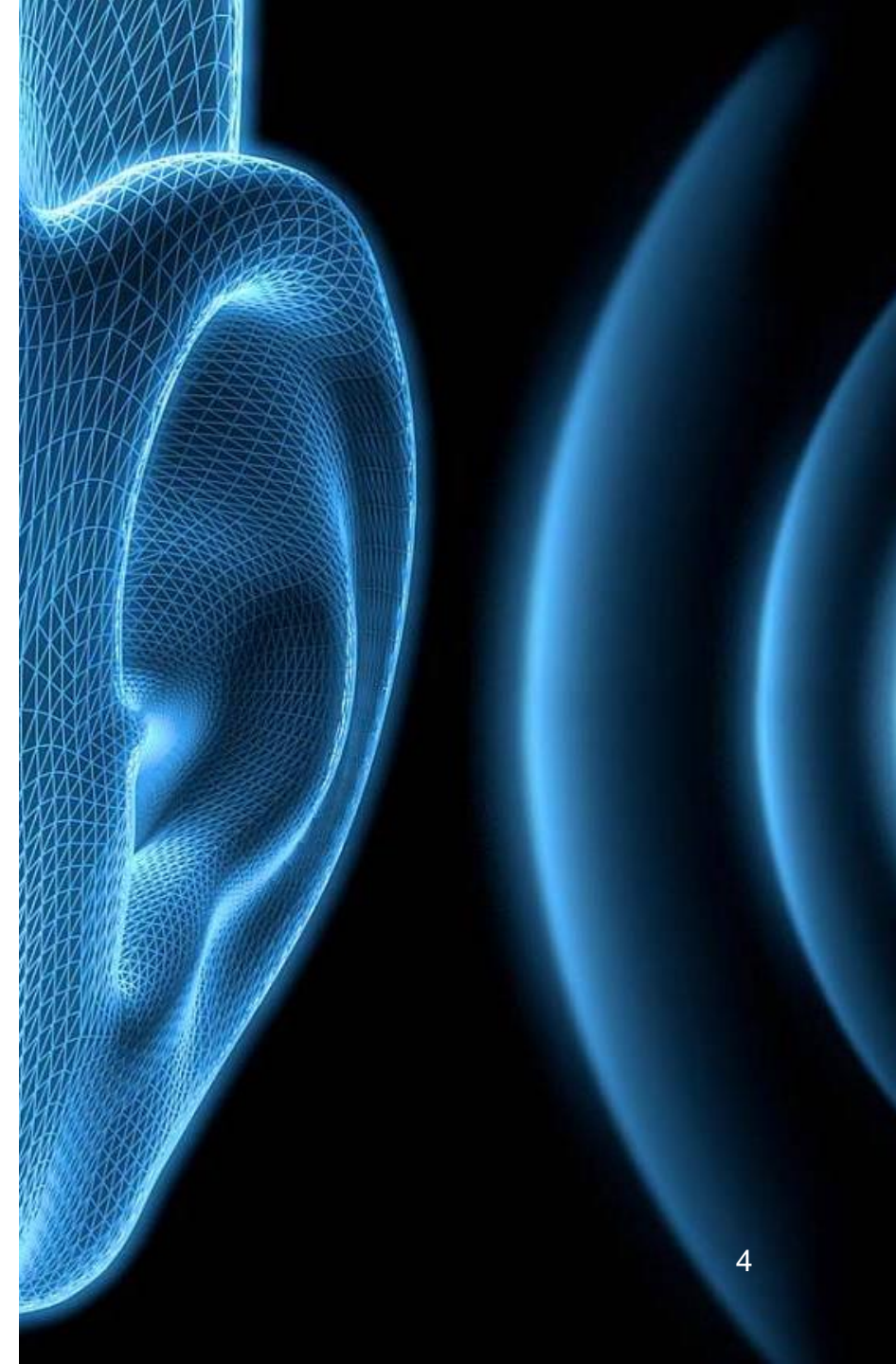
$$\mathbf{w}(k, l) = \frac{\mathbf{R}_n^{-1}(k, l)\mathbf{a}(k, l)}{\mathbf{a}^H(k, l)\mathbf{R}_n^{-1}(k, l)\mathbf{a}(k, l)}$$

Speech Intelligibility in Noise

‘Measure of how comprehensible speech is in given conditions’

Factors that influence speech intelligibility:

- energetic masking
- informational masking
- attention/listening effort
- cognitive abilities/language proficiency
- hearing impairment
- reverberation
- spatial distribution (auditory stream segregation)



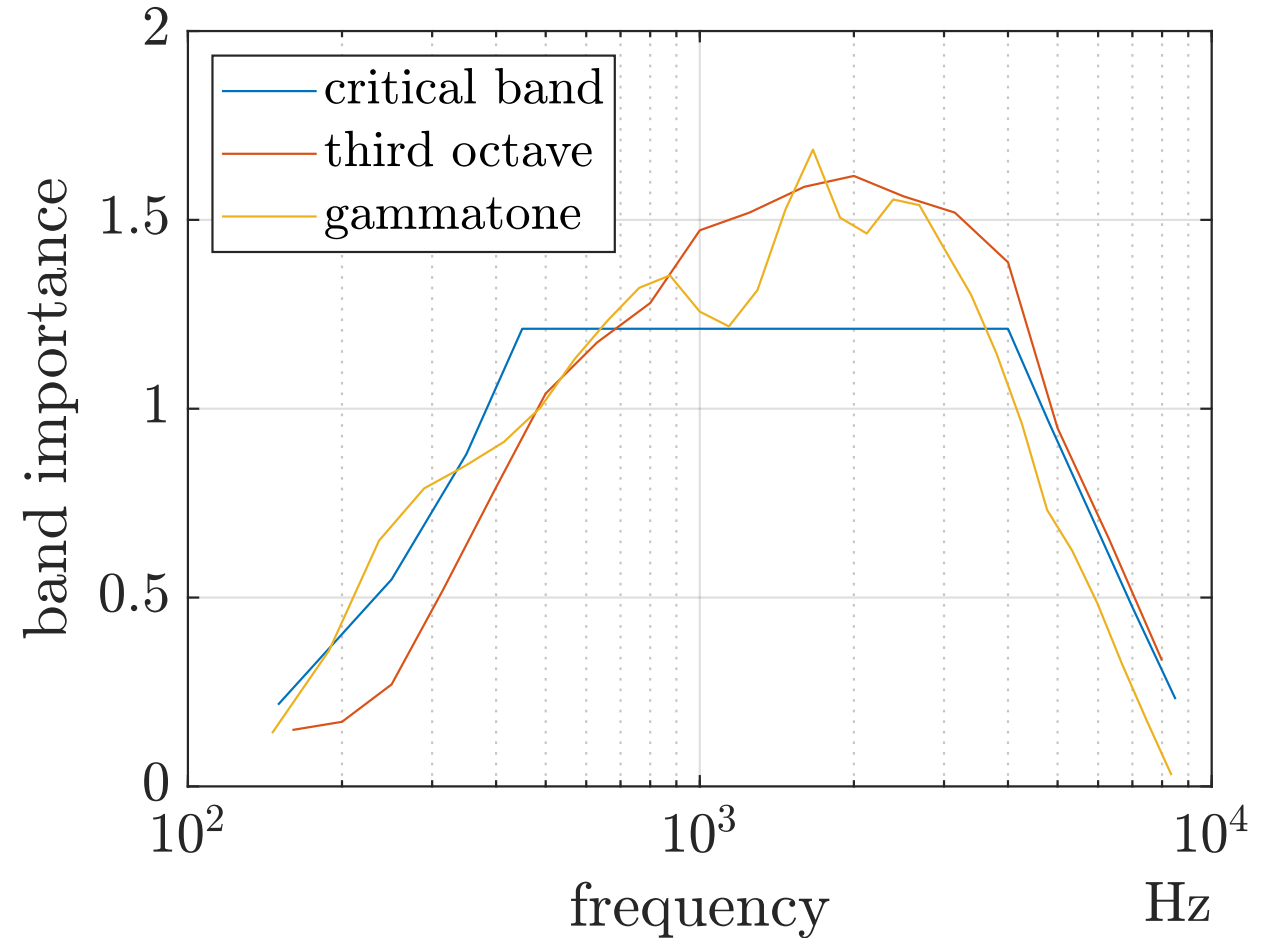


Measuring Intelligibility

- matrix test
- speech reception threshold (SRT)
- audiogram

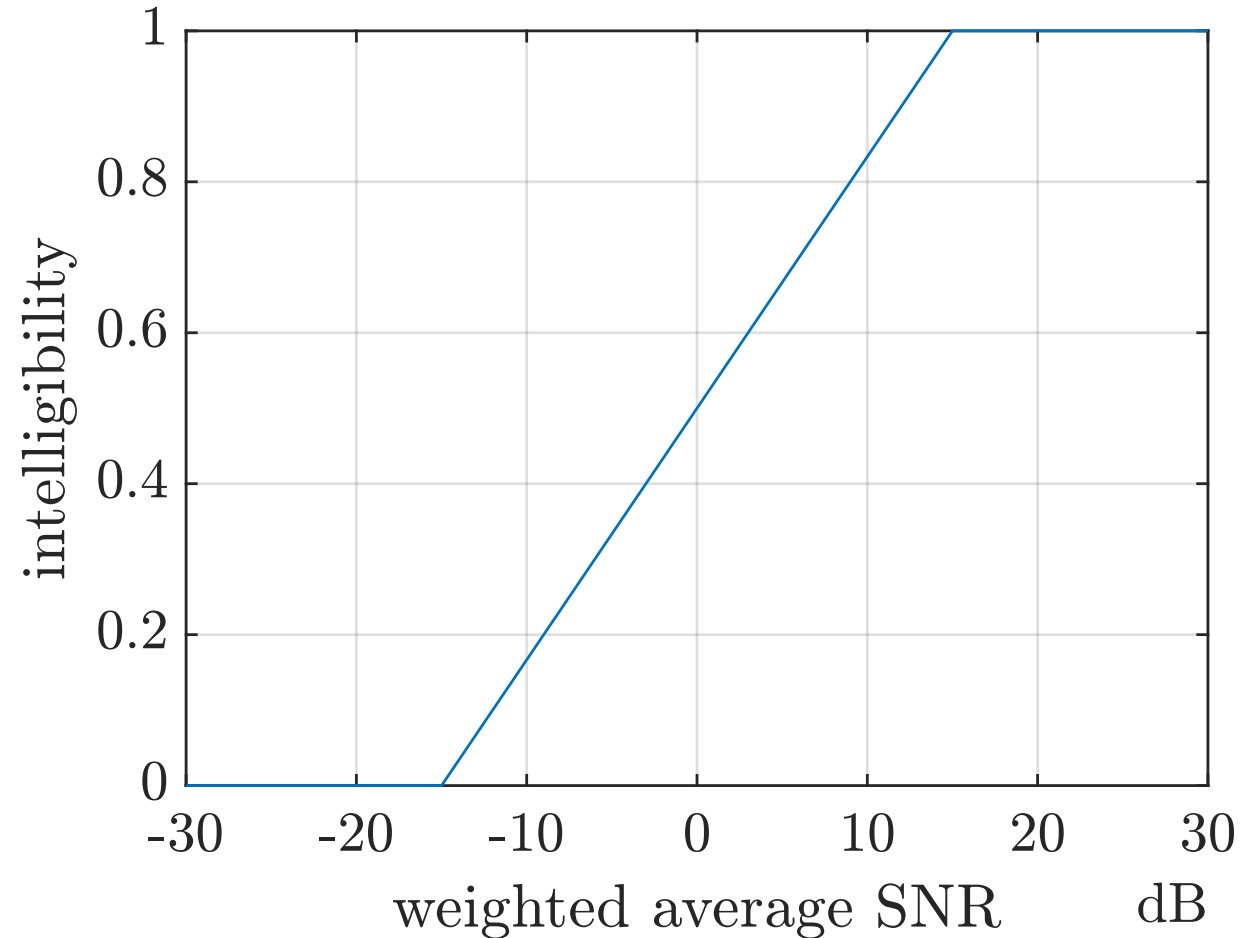
Speech Intelligibility Index (SII)

- objective measure, ANSI standard
- started as articulation index (AI), 1947
- filter speech and noise into frequency bands
- perceptually weighted average of band SNRs
- convert weighted average SNR to intelligibility score



Speech Intelligibility Index (SII)

- objective measure, ANSI standard
- started as articulation index (AI), 1947
- filter speech and noise into frequency bands
- perceptually weighted average of band SNRs
- convert weighted average SNR to intelligibility score
- internal noise to model hearing thresholds
- SRT measure



Other Intelligibility Measures

- speech transmission index (STI)
- coherence speech intelligibility index (CSII)
- short-time objective intelligibility (STOI)
- hearing-aid speech perception index (HASPI)
- data-driven measures

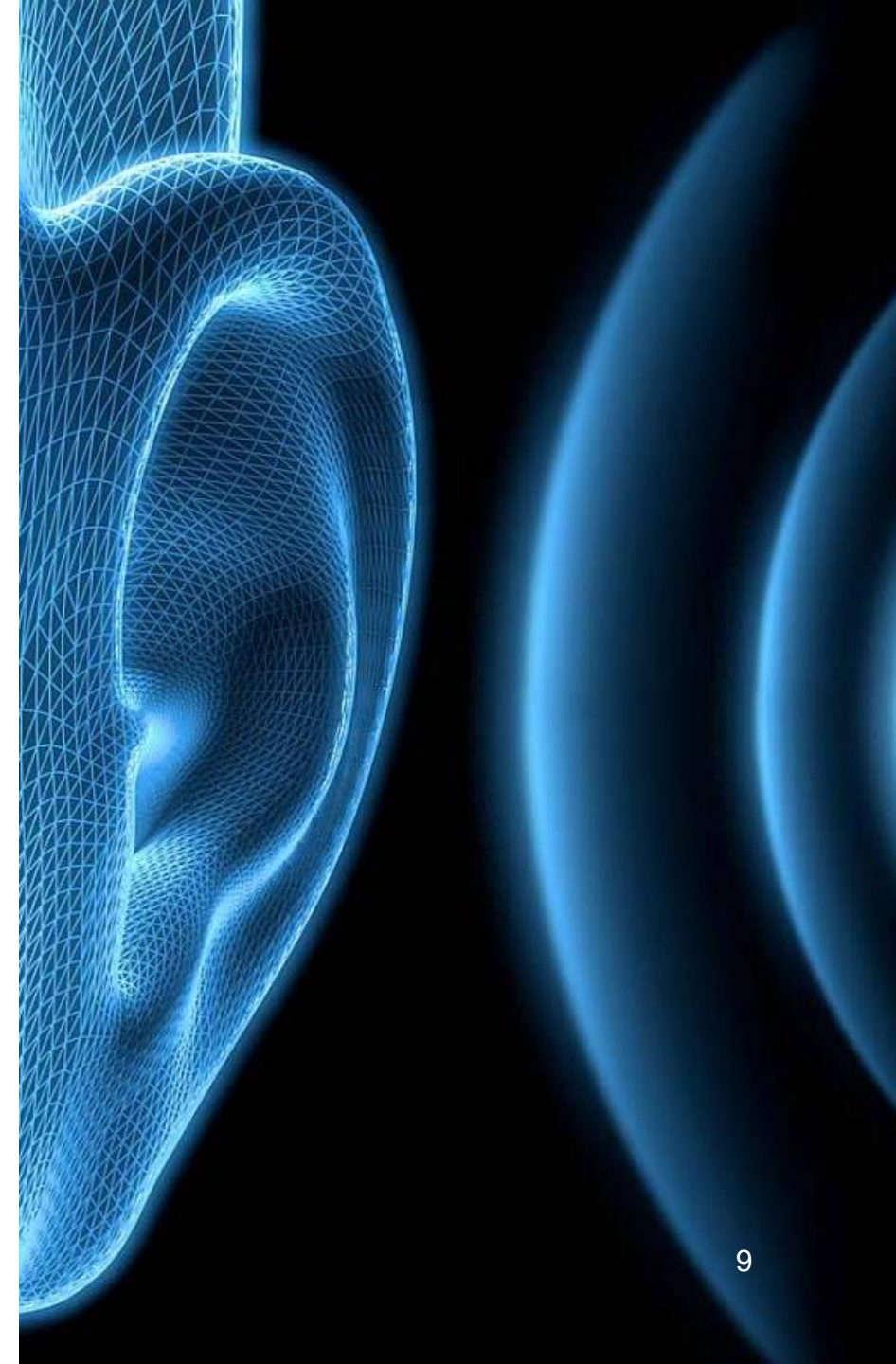


Speech Intelligibility in Noise

‘Measure of how comprehensible speech is in given conditions’

Factors that influence speech intelligibility:

- energetic masking
- informational masking
- attention/listening effort
- cognitive abilities/language proficiency
- hearing impairment
- reverberation
- spatial distribution (auditory stream segregation)



Spatial Release from Masking

two main binaural cues:

- interaural time difference (ITD)
mainly < 1000 Hz
- interaural level difference (ILD)
mainly > 1500 Hz

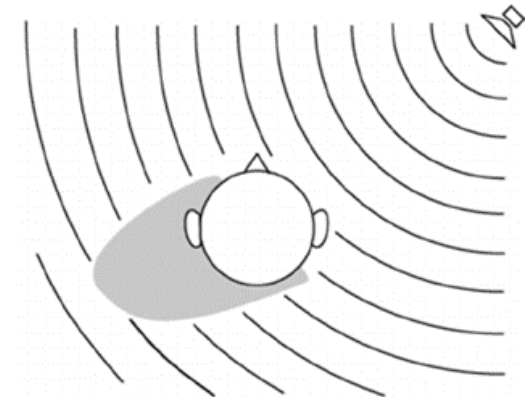
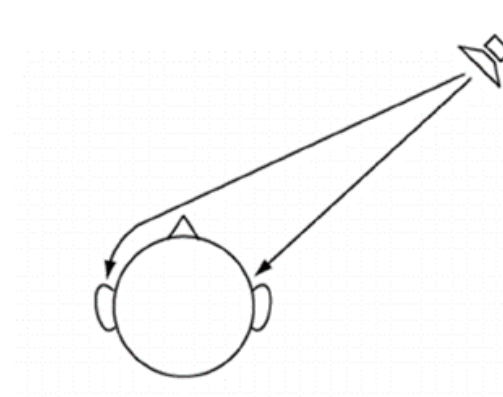
interaural transfer function (ITF)

binaural cue consistency across frequency

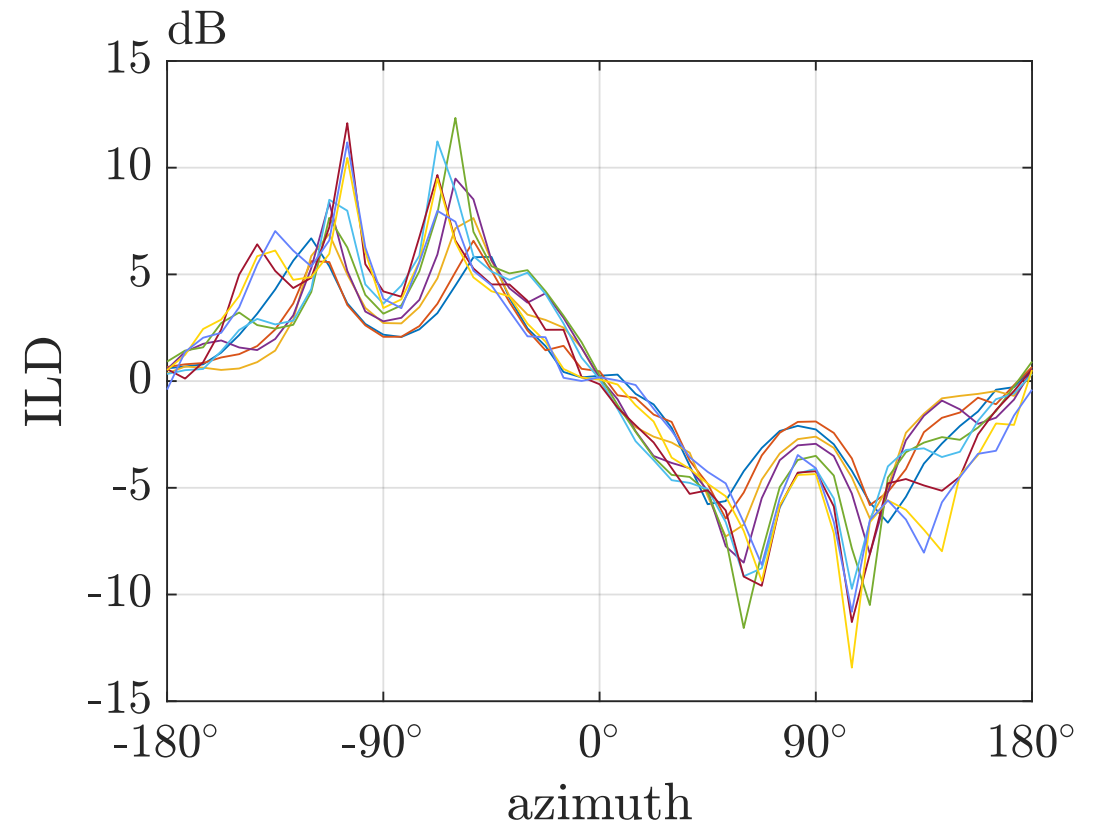
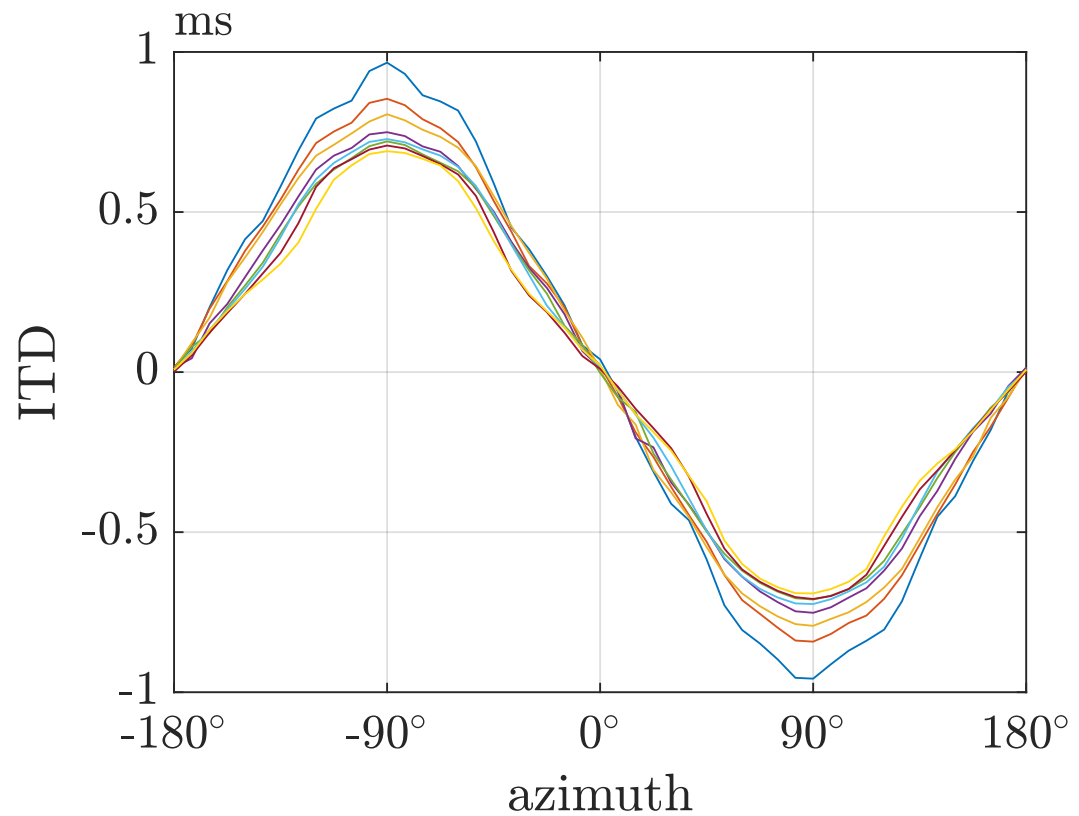
$$\text{ITF} = \frac{a_L}{a_R}$$

$$\text{ITD} = \frac{\angle \text{ITF}}{\omega} = \frac{\angle \text{ITF}}{2\pi f}$$

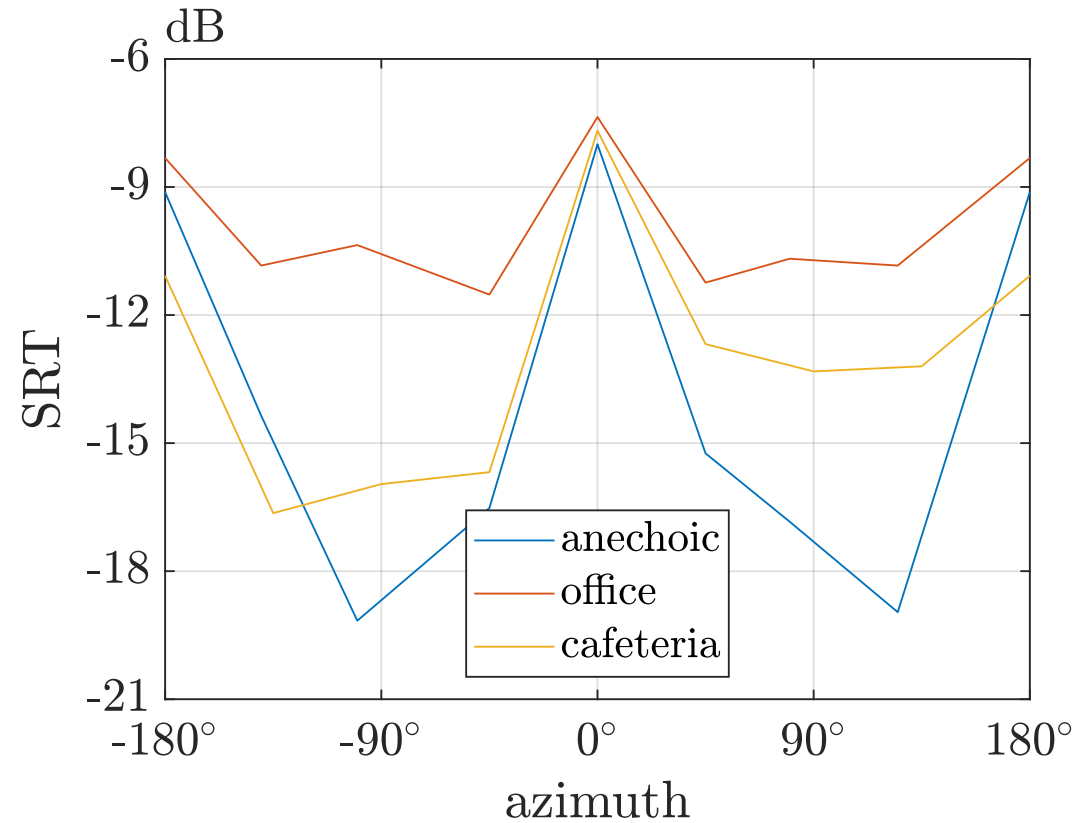
$$\text{ILD} = |\text{ITF}|$$



Spatial Release from Masking



Spatial Release from Masking



Binaural Intelligibility Models

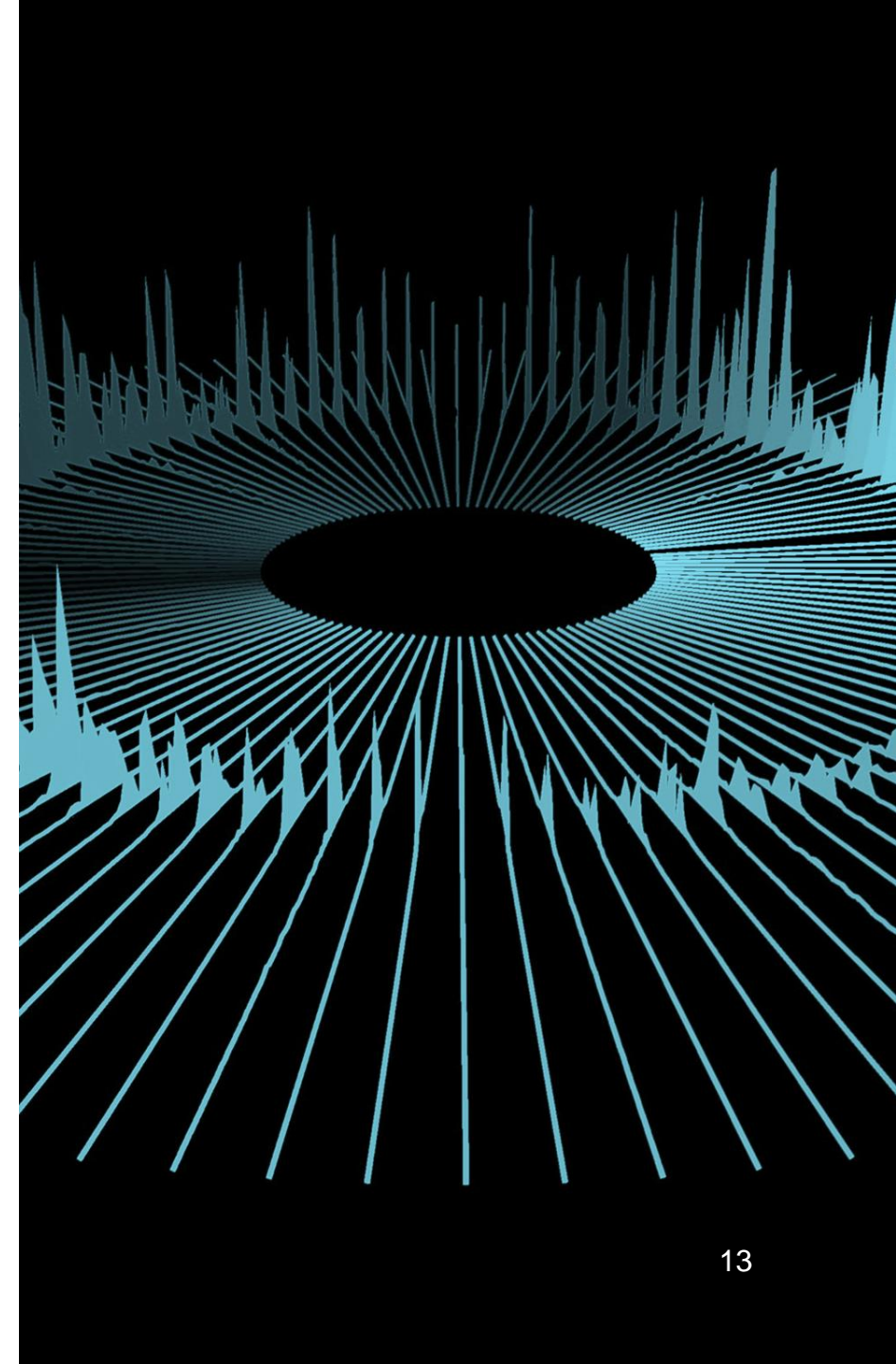
information used:

- SNR (energy)
- modulation
- correlation
- glimpsing

input availability (blindness, intrusiveness)

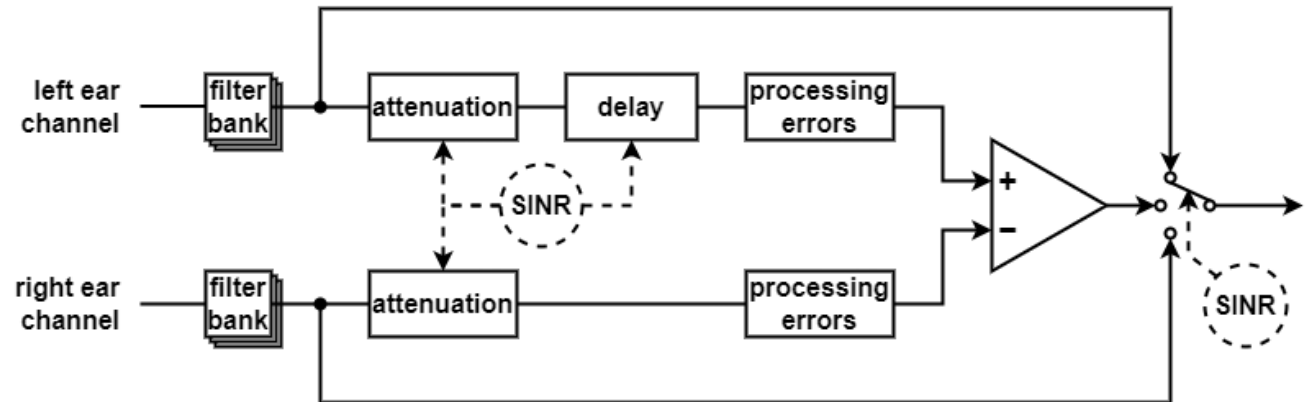
offline/real time

physiological interpretability



Binaural Speech Intelligibility Model

1. gammatone filter bank
2. internal masking noise
3. equalisation–cancellation (Durlach, 1963)
4. artificial processing errors
5. better ear listening
6. monaural intelligibility measure (SII)



Binaural Speech Intelligibility Model

$$y_{EC}(t) = e^{\frac{1}{2}\gamma} y_L(t + \frac{1}{2}\tau) - e^{-\frac{1}{2}\gamma} y_R(t - \frac{1}{2}\tau)$$

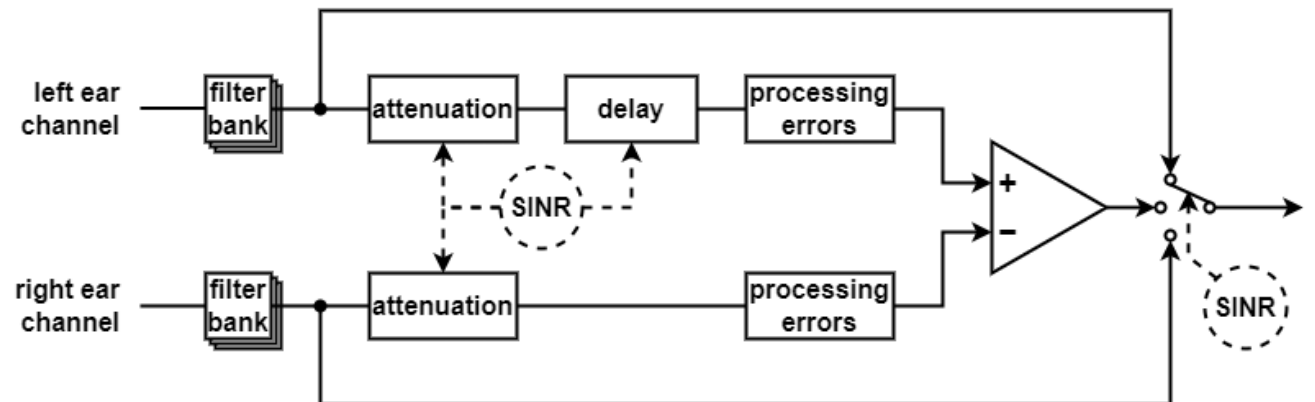
$$y_{EC}(\omega) = e^{\frac{1}{2}(\gamma + j\omega\tau)} y_L(\omega) - e^{-\frac{1}{2}(\gamma + j\omega\tau)} y_R(\omega)$$

$$\mathbf{y} = \begin{pmatrix} y_L \\ y_R \end{pmatrix} \quad \mathbf{v} = \begin{pmatrix} e^{\frac{1}{2}(\gamma - j\omega\tau)} \\ -e^{-\frac{1}{2}(\gamma - j\omega\tau)} \end{pmatrix}$$

$$\text{SNIR}(y_{EC}) = \frac{\mathbf{v}^H \mathbf{P}_s \mathbf{v}}{\mathbf{v}^H \mathbf{P}_n \mathbf{v}}$$

$$y_{EC} = \mathbf{v}^H \mathbf{y}$$

$$\mathbf{P}_y = \int \mathbb{E}(\mathbf{y}\mathbf{y}^H) d\omega$$



Binaural MVDR

binaural signal model

reference transfer function
constraint

What happens to the binaural
cues?

- Target binaural cues are preserved.
- Noise sources virtually move to target!

$$\begin{cases} \mathbf{x} = \mathbf{a}s + \mathbf{n}_0 + \sum_i \mathbf{b}_i n_i \\ y_L = \mathbf{w}_L^H \mathbf{x} \\ y_R = \mathbf{w}_R^H \mathbf{x} \end{cases}$$

$$\begin{array}{ll} \underset{\mathbf{w}_L}{\text{minimise}} & \mathbf{w}_L^H \mathbf{R}_n \mathbf{w}_L \\ \text{subject to} & \mathbf{w}_L^H \mathbf{a} = a_L \end{array} \quad \mathbf{w}_L = \frac{\mathbf{R}_n^{-1} \mathbf{a}}{\mathbf{a}^H \mathbf{R}_n^{-1} \mathbf{a}} a_L^*$$

$$\text{ITF}_{\text{out}}(s) = \frac{\mathbf{w}_L^H \mathbf{a}}{\mathbf{w}_R^H \mathbf{a}} = \frac{a_L}{a_R} = \text{ITF}(s)$$

$$\text{ITF}_{\text{out}}(n_i) = \frac{\mathbf{w}_L^H \mathbf{b}_i}{\mathbf{w}_R^H \mathbf{b}_i} = \frac{a_L}{a_R} = \text{ITF}(s)$$

Joint Binaural LCMV

add binaural cue preservation constraint(s)

couple left and right problems

closed-form solution exists

reduced degrees of freedom for SNR maximisation!

$$\begin{cases} \mathbf{x} = \mathbf{a}s + \mathbf{n}_0 + \sum_i \mathbf{b}_i n_i \\ y_L = \mathbf{w}_L^H \mathbf{x} \\ y_R = \mathbf{w}_R^H \mathbf{x} \end{cases}$$

$$\frac{\mathbf{w}_L^H \mathbf{b}_i}{\mathbf{w}_R^H \mathbf{b}_i} = \frac{b_{iL}}{b_{iR}} \quad \Leftrightarrow \quad \mathbf{w}_L^H \mathbf{b}_i b_{iR} = \mathbf{w}_R^H \mathbf{b}_i b_{iL}$$

$$\begin{aligned} & \underset{\mathbf{w}_L, \mathbf{w}_R}{\text{minimise}} && \mathbf{w}_L^H \mathbf{R}_n \mathbf{w}_L + \mathbf{w}_R^H \mathbf{R}_n \mathbf{w}_R \\ & \text{subject to} && \mathbf{w}_L^H \mathbf{a} = a_L, \\ & && \mathbf{w}_R^H \mathbf{a} = a_R, \\ & && \mathbf{w}_L^H \mathbf{b}_i b_{iR} = \mathbf{w}_R^H \mathbf{b}_i b_{iL} \end{aligned}$$

Perception-Based Beamformer

- include BSIM to account for SRM
- maximise 'perceived' SNR

Binaural MVDR is one of the solutions!

$$\begin{cases} \mathbf{x} = \mathbf{a}s + \mathbf{n}_0 + \sum_i \mathbf{b}_i n_i \\ \mathbf{y} = \mathbf{W}^H \mathbf{x} = \begin{pmatrix} \mathbf{w}_L & \mathbf{w}_R \end{pmatrix}^H \mathbf{x} \\ z = \mathbf{v}^H \mathbf{y} \end{cases}$$

$$\begin{aligned} &\text{maximise}_{\mathbf{v}, \mathbf{W}} \frac{\mathbf{v}^H \mathbf{W}^H \mathbf{a} \mathbf{a}^H \mathbf{W} \mathbf{v}}{\mathbf{v}^H \mathbf{W}^H \mathbf{R}_n \mathbf{W} \mathbf{v}} \\ &\text{subject to} \quad \mathbf{W}^H \mathbf{a} = \mathbf{a}_{\text{ref}} \end{aligned}$$

The logo for TU Delft, featuring a stylized black flame icon above the text 'TU Delft'. The 'TU' is in black, the 'U' is in blue, and 'Delft' is in black.

TU Delft