# End-To-End Deep Learning-Based Adaptation Control for Frequency-Domain Adaptive System Identification
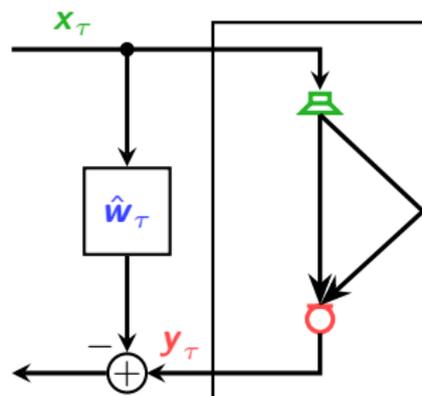
*T. Haubner*, A. Brendel, and W. Kellermann

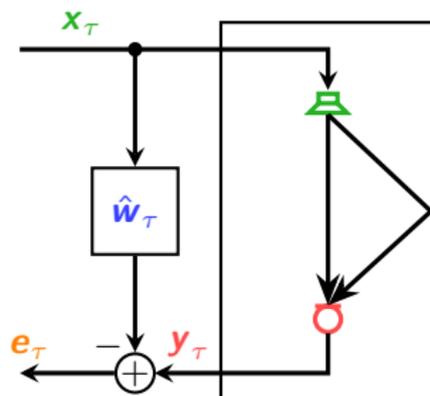**Chair of Multimedia Communications and Signal Processing**

FAU Friedrich-Alexander-Universität Technische Fakultät    icassp 2022 Singapore    LMS

# Motivation: Acoustic Echo Cancellation

▶ **Problem**: Identify acoustic transfer function (ATF) between loudspeaker and microphone signal

# Motivation: Acoustic Echo Cancellation

▶ **Problem**: Identify acoustic transfer function (ATF) between loudspeaker and microphone signal

▶ **Approach:**
  ▶ Minimization of error signal power
  ▶ Iterative update of ATF estimate:
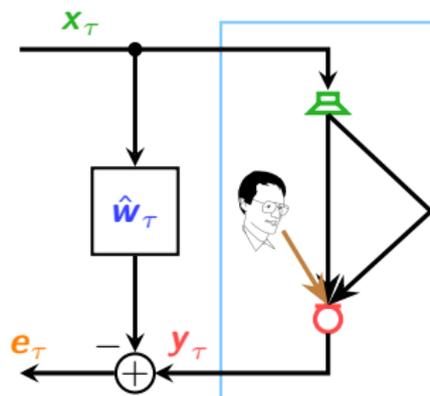    $\hat{\boldsymbol{w}}_\tau = \hat{\boldsymbol{w}}_{\tau-1} + \delta\hat{\boldsymbol{w}}_\tau$

# Motivation: Acoustic Echo Cancellation

▶ **Problem**: Identify acoustic transfer function (ATF)
between loudspeaker and microphone signal

▶ **Approach:**
  - ▶ Minimization of error signal power
  - ▶ Iterative update of ATF estimate:
    $\hat{\boldsymbol{w}}_\tau = \hat{\boldsymbol{w}}_{\tau-1} + \delta\hat{\boldsymbol{w}}_\tau$

▶ **Challenges:**
  - ▶ Interfering signals, e.g., local speech or noise
  - ▶ Time-varying acoustic environments
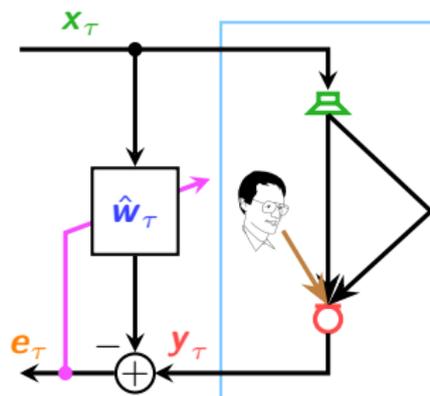
Friedrich-Alexander-Universität
Technische Fakultät

# Motivation: Acoustic Echo Cancellation

▶ **Problem**: Identify acoustic transfer function (ATF) between loudspeaker and microphone signal

▶ **Approach:**
  - ▶ Minimization of error signal power
  - ▶ Iterative update of ATF estimate:
    $\hat{\boldsymbol{w}}_\tau = \hat{\boldsymbol{w}}_{\tau-1} + \delta\hat{\boldsymbol{w}}_\tau$

▶ **Challenges:**
  - ▶ Interfering signals, e.g., local speech or noise
  - ▶ Time-varying acoustic environments



Robust adaptation control for improved convergence rate

Friedrich-Alexander-Universität
Technische Fakultät
Haubner et al.: Deep Learning-Based Adaptation Control
Chair of Multimedia Communications and Signal Processing
May 2022
2 / 20

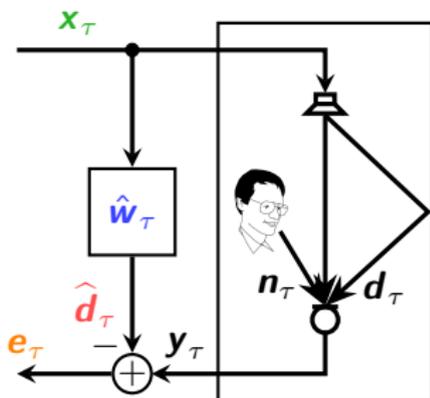# Outline

▶ Acoustic System Identification

▶ Proposed DNN-Based Adaptation Control

▶ Experimental Evaluation

▶ Conclusion

# Iterative ATF Estimation

$$\widehat{\boldsymbol{d}}_\tau = \boldsymbol{Q}(\boldsymbol{x}_\tau \odot \hat{\boldsymbol{w}}_{\tau-1})$$

1. Estimate echo $\widehat{\boldsymbol{d}}_\tau$ by linear convolution of
   - ▶ loudspeaker signal block $\boldsymbol{x}_\tau$
   - ▶ with ATF estimate $\hat{\boldsymbol{w}}_{\tau-1}$.
   - ▶ and linear convolution constraint matrix $\boldsymbol{Q}$



Haubner et al.: Deep Learning-Based Adaptation Control
Chair of Multimedia Communications and Signal Processing
May 2022
4 / 20

# Iterative ATF Estimation

$$\boldsymbol{e}_\tau = \boldsymbol{y}_\tau - \widehat{\boldsymbol{d}}_\tau = \boldsymbol{y}_\tau - \boldsymbol{Q}(\boldsymbol{x}_\tau \odot \hat{\boldsymbol{w}}_{\tau-1})$$

1. Estimate echo $\widehat{\boldsymbol{d}}_\tau$ by linear convolution of
   - loudspeaker signal block $\boldsymbol{x}_\tau$
   - with ATF estimate $\hat{\boldsymbol{w}}_{\tau-1}$.
   - and linear convolution constraint matrix $\boldsymbol{Q}$
2. Compute error signal block $\boldsymbol{e}_\tau$

# Iterative ATF Estimation

$$\boldsymbol{e}_\tau = \boldsymbol{y}_\tau - \widehat{\boldsymbol{d}}_\tau = \boldsymbol{y}_\tau - \boldsymbol{Q}(\boldsymbol{x}_\tau \odot \hat{\boldsymbol{w}}_{\tau-1})$$
$$\hat{\boldsymbol{w}}_\tau = \hat{\boldsymbol{w}}_{\tau-1} + \boldsymbol{G} \quad (\boldsymbol{x}_\tau^* \odot \boldsymbol{e}_\tau)$$

1. Estimate echo $\widehat{\boldsymbol{d}}_\tau$ by linear convolution of
   - ▶ loudspeaker signal block $\boldsymbol{x}_\tau$
   - ▶ with ATF estimate $\hat{\boldsymbol{w}}_{\tau-1}$.
   - ▶ and linear convolution constraint matrix $\boldsymbol{Q}$
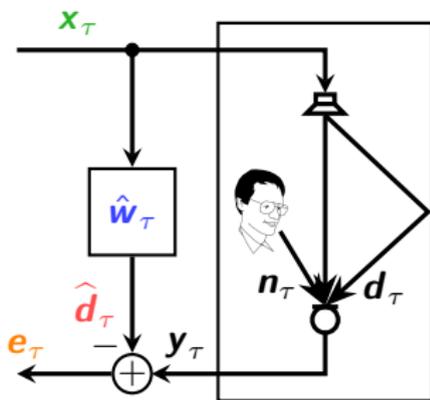2. Compute error signal block $\boldsymbol{e}_\tau$
3. Compute stochastic gradient $\boldsymbol{G}\,(\boldsymbol{x}_\tau^* \odot \boldsymbol{e}_\tau)$
   with FIR filter projection matrix $\boldsymbol{G}$
4. Perform gradient descent



Haubner et al.: Deep Learning-Based Adaptation Control
Chair of Multimedia Communications and Signal Processing

May 2022
4 / 20

Friedrich-Alexander-Universität
Technische Fakultät

# Adaptation Control

$$e_\tau = y_\tau - \widehat{d}_\tau = n_\tau + \left( d_\tau - \widehat{d}_\tau \right)$$

$$\hat{w}_\tau = \hat{w}_{\tau-1} + G \quad \left( x_\tau^* \odot e_\tau \right)$$

Large error powers $\|e_\tau\|^2$ could result from

▶ system mismatch, i.e., imprecise echo estimates $\widehat{d}_\tau$
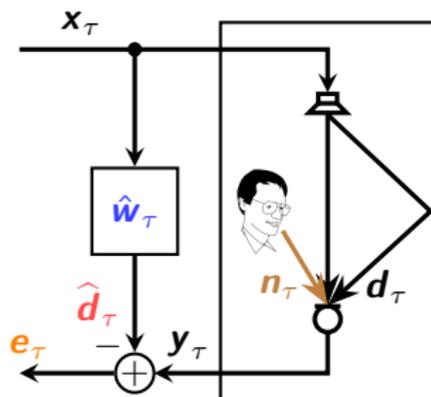
▶ interfering signals $n_\tau$

# Adaptation Control

$$e_\tau = y_\tau - \widehat{d}_\tau = n_\tau + \left( d_\tau - \widehat{d}_\tau \right)$$

$$\hat{w}_\tau = \hat{w}_{\tau-1} + G \quad \left( x_\tau^* \odot e_\tau \right)$$

Large error powers $||e_\tau||^2$ could result from

▶ system mismatch, i.e., imprecise echo
  estimates $\widehat{d}_\tau \Rightarrow$ Update filter coefficients

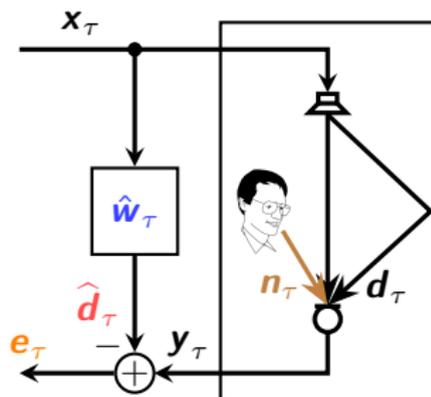▶ interfering signals $n_\tau \Rightarrow$ Stall filter adaptation

# Adaptation Control

$$e_\tau = y_\tau - \widehat{d}_\tau = n_\tau + \left( d_\tau - \widehat{d}_\tau \right)$$
$$\hat{w}_\tau = \hat{w}_{\tau-1} + G\Lambda_\tau (x_\tau^* \odot e_\tau)$$

Large error powers $||e_\tau||^2$ could result from

▶ system mismatch, i.e., imprecise echo estimates $\widehat{d}_\tau$ ⇒ Update filter coefficients

▶ interfering signals $n_\tau$ ⇒ Stall filter adaptation

**Solution**: Control adaptation by time- and frequency-dependent step-sizes $[\Lambda_\tau]_{ff}$



Haubner et al.: Deep Learning-Based Adaptation Control
Chair of Multimedia Communications and Signal Processing     May 2022
5 / 20

Friedrich-Alexander-Universität
Technische Fakultät

# Traditional Model-Based Adaptation Control

$$\mathbf{\Lambda}_\tau^{\mathsf{MB}} = f_{\mathsf{MB}}\left(\mathbf{\psi}_\tau^{\mathsf{XX}}, \mathbf{\psi}_\tau^{\mathsf{NN}}, \mathbf{\psi}_\tau^{\Delta\mathsf{W}\Delta\mathsf{W}}, \dots\right)$$

Compute step-size matrix $\mathbf{\Lambda}_\tau^{\mathsf{MB}}$ as a function of

▶ loudspeaker PSD $\mathbf{\psi}_\tau^{\mathsf{XX}}$

▶ interference PSD $\mathbf{\psi}_\tau^{\mathsf{NN}}$

▶ filter estimation uncertainty $\mathbf{\psi}_\tau^{\Delta\mathsf{W}\Delta\mathsf{W}}$

▶ . . .

**Prominent examples**:
FDAF, DFT-domain Kalman filter, . . .

# Traditional Model-Based Adaptation Control

$$\mathbf{\Lambda}_\tau^{\mathrm{MB}} = f_{\mathrm{MB}}\left(\mathbf{\psi}_\tau^{\mathrm{XX}}, \mathbf{\psi}_\tau^{\mathrm{NN}}, \mathbf{\psi}_\tau^{\Delta\mathrm{W}\Delta\mathrm{W}}, \dots\right)$$

**Challenges**

▶ Estimation of signal statistics of unknown quantities, e.g., $\mathbf{\psi}_\tau^{\mathrm{NN}}$

▶ Mismatch of assumed model properties

FAU Friedrich-Alexander-Universität
Technische Fakultät

Haubner et al.: Deep Learning-Based Adaptation Control
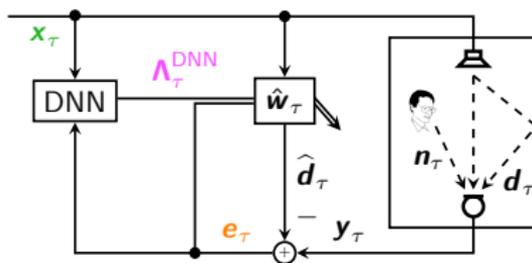Chair of Multimedia Communications and Signal Processing

May 2022
6 / 20

LMS

# Outline

▶ Acoustic System Identification

▶ Proposed DNN-Based Adaptation Control

▶ Experimental Evaluation

▶ Conclusion

FAU Friedrich-Alexander-Universität
Technische Fakultät    Haubner et al.:   Deep Learning-Based Adaptation Control    May 2022
Chair of Multimedia Communications and Signal Processing    7 / 20    LMS

# Proposed General Concept

$$\boldsymbol{\Lambda}_{\tau}^{\text{DNN}} = f_{\text{DNN}}\left(\boldsymbol{x}_1, \boldsymbol{e}_1, \ldots, \boldsymbol{x}_{\tau}, \boldsymbol{e}_{\tau}; \boldsymbol{\theta}\right)$$

▶ Learn mapping $f_{\text{DNN}}$ of observed signal sequences to step-size matrices

▶ DNN parameters $\boldsymbol{\theta}$ are learned from training data

FAU Friedrich-Alexander-Universität
Technische Fakultät

Haubner et al.:   Deep Learning-Based Adaptation Control
Chair of Multimedia Communications and Signal Processing

May 2022
8 / 20

LMS

# Proposed Step-Size Structure

$$\left[\mathbf{\Lambda}_\tau^{\mathsf{DNN}}\right]_{ff} = m_{f,\tau}^\mu$$

DNN provides raw step-size $m_{f,\tau}^\mu \in [0, \mu_{\mathsf{MAX}}]$

# Proposed Step-Size Structure

$$\left[\mathbf{\Lambda}_\tau^{\mathrm{DNN}}\right]_{ff} = m_{f,\tau}^\mu$$

**Challenges:** DNN needs to model

▶ large numerical range due to non-whiteness of loudspeaker signals

DNN provides raw step-size $m_{f,\tau}^\mu \in [0, \mu_{\mathsf{MAX}}]$

# Proposed Step-Size Structure

$$\left[ \boldsymbol{\Lambda}_{\tau}^{\mathsf{DNN}} \right]_{ff} = \frac{m_{f,\tau}^{\mu}}{\hat{\boldsymbol{\psi}}_{f,\tau}^{\mathsf{XX}}}$$

**Challenges:** DNN needs to model

▶ large numerical range due to non-whiteness of loudspeaker signals

Normalize raw step-size $m_{f,\tau}^{\mu}$ by **loudspeaker PSD** estimate $\hat{\boldsymbol{\psi}}_{f,\tau}^{\mathsf{XX}}$

## Proposed Step-Size Structure

$$\left[ \boldsymbol{\Lambda}_\tau^{\mathsf{DNN}} \right]_{ff} = \frac{m_{f,\tau}^\mu}{\hat{\boldsymbol{\Psi}}_{f,\tau}^{\mathsf{XX}}}$$

**Challenges:** DNN needs to model

▶ large numerical range due to non-whiteness of loudspeaker signals

▶ rapid changes due to non-stationarity of interfering signals

Friedrich-Alexander-Universität
Technische Fakultät
Haubner et al.: Deep Learning-Based Adaptation Control
Chair of Multimedia Communications and Signal Processing
May 2022
9 / 20
LMS

# Proposed Step-Size Structure

$$\left[ \mathbf{\Lambda}_\tau^{\mathsf{DNN}} \right]_{ff} = \frac{m_{f,\tau}^\mu}{\hat{\Psi}_{f,\tau}^{\mathsf{XX}} + \hat{\Psi}_{f,\tau}^{\mathsf{II}}}$$

**Challenges:** DNN needs to model

▶ large numerical range due to non-whiteness of loudspeaker signals

▶ rapid changes due to non-stationarity of interfering signals

**Traditional Approach**: Normalization by interference PSD estimate $\hat{\Psi}_{f,\tau}^{\mathsf{II}}$

## Proposed Step-Size Structure

$$\left[ \mathbf{\Lambda}_\tau^{\mathsf{DNN}} \right]_{ff} = \frac{m_{f,\tau}^\mu}{\hat{\Psi}_{f,\tau}^{\mathsf{XX}} + \frac{M}{R} |m_{f,\tau}^e \left[ \mathbf{e}_\tau \right]_f |^2}$$

**Challenges:** DNN needs to model

▶ large numerical range due to non-whiteness of loudspeaker signals

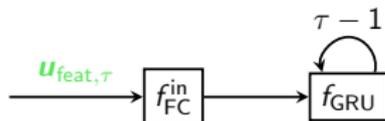▶ rapid changes due to non-stationarity of interfering signals

Normalization by **masked error power** $|m_{f,\tau}^e \left[ \mathbf{e}_\tau \right]_f |^2$

**FAU** Friedrich-Alexander-Universität
Technische Fakultät

Haubner et al.:  Deep Learning-Based Adaptation Control          May 2022
Chair of Multimedia Communications and Signal Processing      9 / 20

**LMS**

# Proposed Step-Size Structure

$$\left[\mathbf{\Lambda}_\tau^{\mathsf{DNN}}\right]_{ff} = \frac{m_{f,\tau}^\mu}{\hat{\Psi}_{f,\tau}^{\mathsf{XX}} + \frac{M}{R}|m_{f,\tau}^e\,[\boldsymbol{e}_\tau]_f\,|^2}$$

**Challenges:** DNN needs to model

▶ large numerical range due to non-whiteness of loudspeaker signals

▶ rapid changes due to non-stationarity of interfering signals

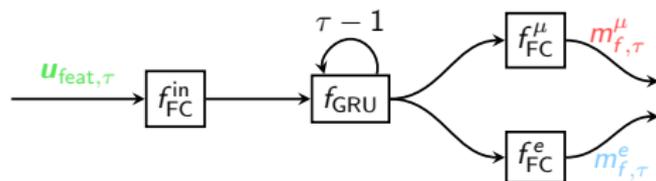Exploitation of **domain knowledge** from traditional adaptation control

**FAU** Friedrich-Alexander-Universität
Technische Fakultät
Haubner et al.: Deep Learning-Based Adaptation Control
Chair of Multimedia Communications and Signal Processing
May 2022
9 / 20
**LMS**

# DNN Architecture

$$\xrightarrow{\boldsymbol{u}_{\text{feat},\tau}} \boxed{f^{\text{in}}_{\text{FC}}}$$
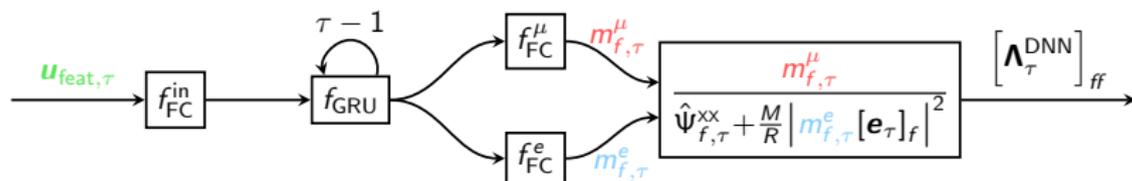
- ▶ Feature vector $\boldsymbol{u}_{\text{feat},\tau}$
  - ▶ Log. loudspeaker power spectrum
  - ▶ Log. error power spectrum
- ▶ Extract temporal information by GRU layers
- ▶ DNN provides masks $m^{\mu}_{f,\tau} \in [0, \mu_{\text{MAX}}]$ and $m^{e}_{f,\tau} \in [0, 1]$
- ▶ Time- and frequency-dependent step-sizes $[\boldsymbol{\Lambda}^{\text{DNN}}_{\tau}]_{ff}$
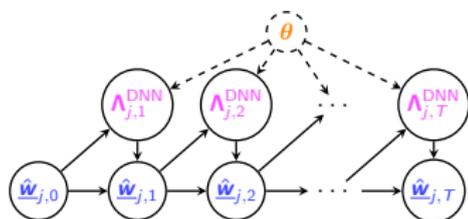
# DNN Architecture



- ▶ Feature vector $\boldsymbol{u}_{\text{feat},\tau}$
    - ▶ Log. loudspeaker power spectrum
    - ▶ Log. error power spectrum
- ▶ Extract temporal information by GRU layers
- ▶ DNN provides masks $m^{\mu}_{f,\tau} \in [0, \mu_{\text{MAX}}]$ and $m^{e}_{f,\tau} \in [0, 1]$
- ▶ Time- and frequency-dependent step-sizes $[\boldsymbol{\Lambda}^{\text{DNN}}_{\tau}]_{ff}$

Haubner et al.:  Deep Learning-Based Adaptation Control     May 2022
Chair of Multimedia Communications and Signal Processing     10 / 20

FAU Friedrich-Alexander-Universität
Technische Fakultät

LMS

# DNN Architecture



- ▶ Feature vector $\boldsymbol{u}_{\text{feat},\tau}$
  - ▶ Log. loudspeaker power spectrum
  - ▶ Log. error power spectrum
- ▶ Extract temporal information by GRU layers
- ▶ DNN provides masks $m_{f,\tau}^{\mu} \in [0, \mu_{\text{MAX}}]$ and $m_{f,\tau}^{e} \in [0, 1]$
- ▶ Time- and frequency-dependent step-sizes $[\boldsymbol{\Lambda}_{\tau}^{\text{DNN}}]_{ff}$

# DNN Architecture



- Feature vector $\boldsymbol{u}_{\text{feat},\tau}$
  - Log. loudspeaker power spectrum
  - Log. error power spectrum
- Extract temporal information by GRU layers
- DNN provides masks $m_{f,\tau}^{\mu} \in [0, \mu_{\text{MAX}}]$ and $m_{f,\tau}^{e} \in [0,1]$
- Time- and frequency-dependent step-sizes $[\boldsymbol{\Lambda}_{\tau}^{\text{DNN}}]_{ff}$

# DNN Training

**Challenges:**

▶ Choice of optimum target step-sizes

▶ Cost function design

**FAU** Friedrich-Alexander-Universität
Technische Fakultät

Haubner et al.:  Deep Learning-Based Adaptation Control    May 2022
Chair of Multimedia Communications and Signal Processing    11 / 20

**LMS**

# Proposed Solution

> End-to-end training of DNN parameters $\boldsymbol{\theta}$ w.r.t.
> achieved system identification performance

Cost function:

$$\mathcal{J}(\boldsymbol{\theta}) = \frac{1}{TJ} \sum_{j=1}^{J} \sum_{\tau=1}^{T} 10 \log_{10} \left( \Upsilon_{j,\tau} \right)$$

with normalized system distance

$$\Upsilon_{j,\tau} = \frac{||\underline{\boldsymbol{w}}_{j,\tau} - \underline{\hat{\boldsymbol{w}}}_{j,\tau}||_2^2}{||\underline{\boldsymbol{w}}_{j,\tau}||_2^2}$$

$j$: sequence index, $\tau$: block index

Friedrich-Alexander-Universität
Technische Fakultät

# Outline

▶ Acoustic System Identification

▶ Proposed DNN-Based Adaptation Control

▶ **Experimental Evaluation**

▶ Conclusion

FAU Friedrich-Alexander-Universität
Technische Fakultät          Haubner et al.:   Deep Learning-Based Adaptation Control          May 2022
          Chair of Multimedia Communications and Signal Processing          13 / 20          LMS

# Experimental Evaluation: AEC Application

▶ Loudspeaker signal: 143 different speakers

▶ Acoustic environment: 201 measured AIRs

    ▶ Sampling frequency: $f_s = 16$ kHz

    ▶ Reverberation time $T_{60} \in [120\text{ms}, 780\text{ms}]$

    ▶ Acoustic scene change in the interval [7.2s, 8.8s]

▶ Interfering signal

    ▶ 145 different speakers
    (echo-to-near-end power ratio $\in [-10\text{dB}, 10\text{dB}]$)

    ▶ White Gaussian noise
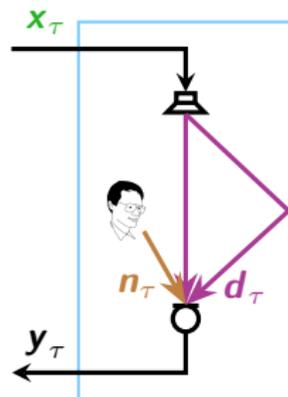    (echo-to-noise power ratio $\in [25\text{dB}, 35\text{dB}]$)

Friedrich-Alexander-Universität
Technische Fakultät

Haubner et al.: Deep Learning-Based Adaptation Control          May 2022
Chair of Multimedia Communications and Signal Processing          14 / 20

LMS

# Experimental Evaluation: AEC Application
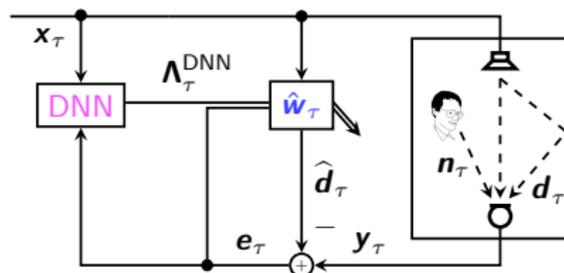
► Loudspeaker signal: 143 different speakers

► Acoustic environment: 201 measured AIRs

  ► Sampling frequency: $f_s = 16$ kHz

  ► Reverberation time $T_{60} \in [120\text{ms}, 780\text{ms}]$

  ► Acoustic scene change in the interval [7.2s, 8.8s]

► Interfering signal

  ► 145 different speakers
  (echo-to-near-end power ratio ∈ [−10dB, 10dB])

  ► White Gaussian noise
  (echo-to-noise power ratio ∈ [25dB, 35dB])

FAU Friedrich-Alexander-Universität
Technische Fakultät

Haubner et al.: Deep Learning-Based Adaptation Control        May 2022
Chair of Multimedia Communications and Signal Processing        14 / 20

LMS

# Experimental Evaluation: AEC Application

▶ Loudspeaker signal: 143 different speakers

▶ Acoustic environment: 201 measured AIRs
  ▶ Sampling frequency: $f_s = 16$ kHz
  ▶ Reverberation time $T_{60} \in$ [120ms, 780ms]
  ▶ Acoustic scene change in the interval [7.2s, 8.8s]

▶ Interfering signal
  ▶ 145 different speakers
    (echo-to-near-end power ratio $\in$ [−10dB, 10dB])
  ▶ White Gaussian noise
    (echo-to-noise power ratio $\in$ [25dB, 35dB])

# Experimental Evaluation: AEC Application

▶ Loudspeaker signal: 143 different speakers

▶ Acoustic environment: 201 measured AIRs

    ▶ Sampling frequency: $f_s = 16$ kHz

    ▶ Reverberation time $T_{60} \in [120\text{ms}, 780\text{ms}]$

    ▶ Acoustic scene change in the interval [7.2s, 8.8s]

▶ Interfering signal

    ▶ 145 different speakers
    (echo-to-near-end power ratio $\in [-10\text{dB}, 10\text{dB}]$)

    ▶ White Gaussian noise
    (echo-to-noise power ratio $\in [25\text{dB}, 35\text{dB}]$)



Separation into disjoint training and testing data sets

# Algorithmic Settings

- **FIR filter** length: 2048 samples

- Frame shift: 1024 samples

- #DNN parameters: 2.4 million
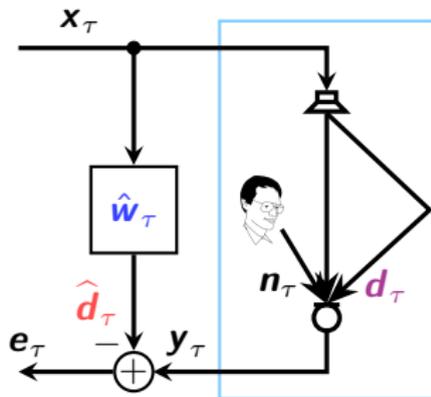
- Training optimizer: *Adam*

# Performance Measures

▶ **Normalized system distance** (the lower the better)

$$\Upsilon_{\text{ZP},\tau} = 10 \log_{10} \frac{\|\tilde{\underline{w}}_\tau - V\hat{\underline{w}}_\tau\|_2^2}{\|\tilde{\underline{w}}_\tau\|_2^2}$$
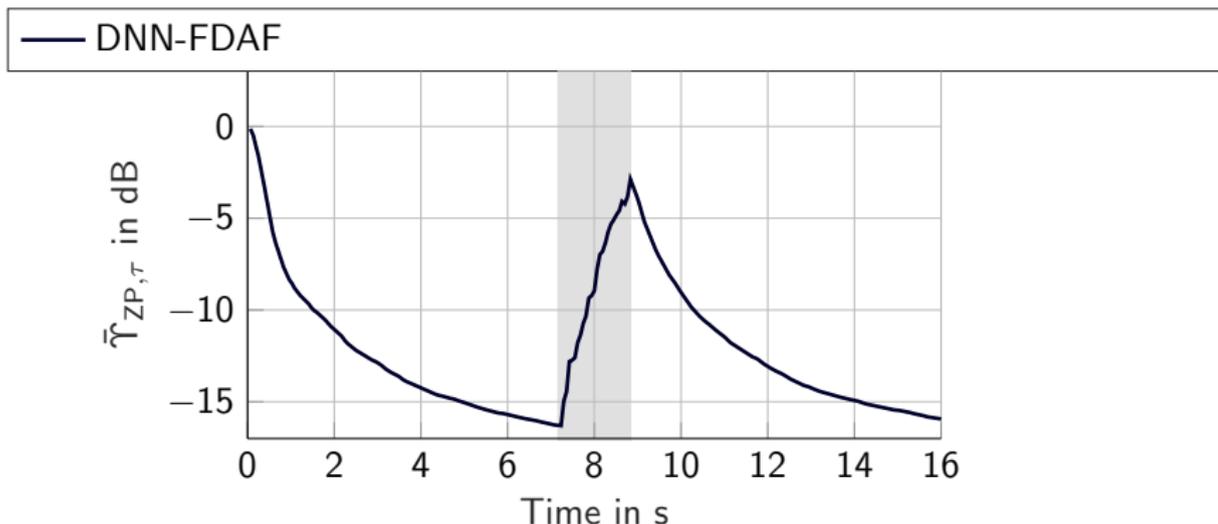
with $V$ being a zero-padding matrix

▶ **ERLE** (the higher the better)

$$\mathcal{E}_\tau = 10 \log_{10} \frac{\mathbb{E}\left[\|\underline{d}_\tau\|_2^2\right]}{\mathbb{E}\left[\|\underline{d}_\tau - \hat{\underline{d}}_\tau\|_2^2\right]}$$

# Performance Measures

▶ **Normalized system distance** (the lower the better)

$$\Upsilon_{\text{ZP},\tau} = 10 \log_{10} \frac{||\tilde{\underline{w}}_\tau - V\hat{\underline{w}}_\tau||_2^2}{||\tilde{\underline{w}}_\tau||_2^2}$$

with $V$ being a zero-padding matrix

▶ **ERLE** (the higher the better)

$$\mathcal{E}_\tau = 10 \log_{10} \frac{\mathbb{E}\left[||\underline{d}_\tau||_2^2\right]}{\mathbb{E}\left[||\underline{d}_\tau - \hat{\underline{d}}_\tau||_2^2\right]}$$

Arithmetic average of **100** experiments with randomly-selected loudspeaker and interfering signals, AIRs and transition times.
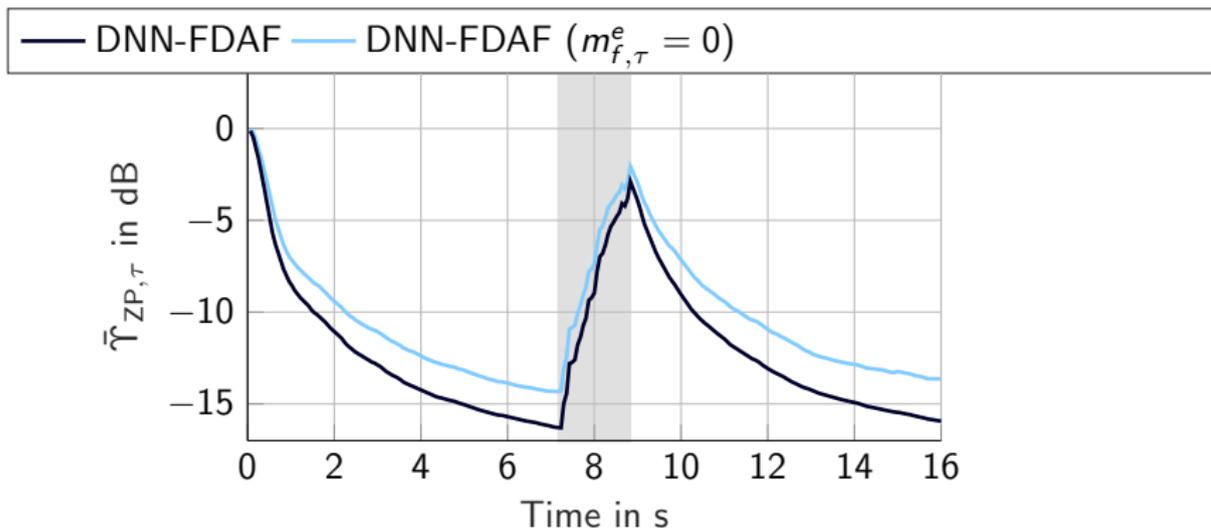
## Analysis of Proposed DNN-FDAF



$$\left[\mathbf{\Lambda}_\tau^{\mathrm{DNN}}\right]_{ff} = \frac{m_{f,\tau}^\mu}{\hat{\Psi}_{f,\tau}^{\mathrm{XX}} + \frac{M}{R}\left|m_{f,\tau}^e\left[\boldsymbol{e}_\tau\right]_f\right|^2}$$

**Proposed** step-size control

## Analysis of Proposed DNN-FDAF



$$\left[ \mathbf{\Lambda}_\tau^{\mathsf{DNN}} \right]_{ff} = \frac{m_{f,\tau}^\mu}{\hat{\psi}_{f,\tau}^{\mathsf{XX}} + \; 0}$$

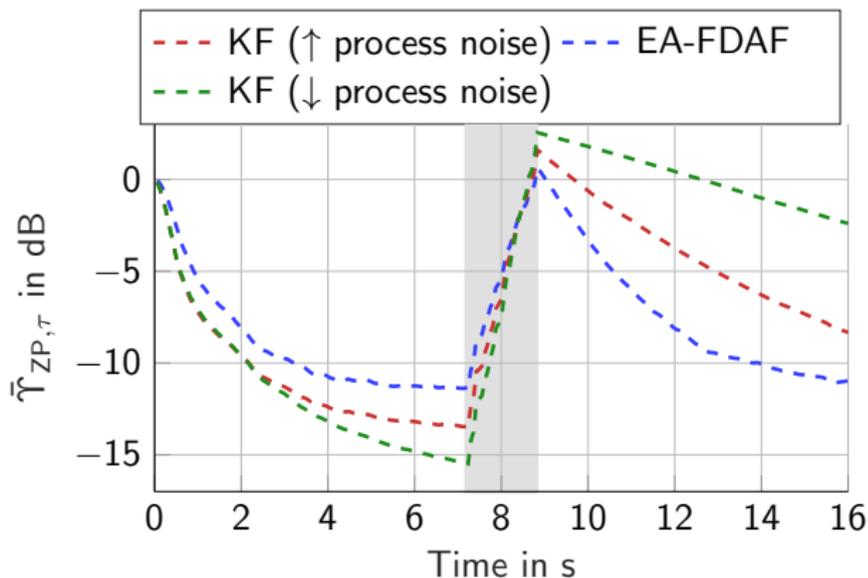Discarding error power normalization

## Analysis of Proposed DNN-FDAF



$$\left[\mathbf{\Lambda}_\tau^{\text{DNN}}\right]_{ff} = \frac{0.5}{\hat{\Psi}_{f,\tau}^{XX} + \frac{M}{R}\left|m_{f,\tau}^e\left[\mathbf{e}_\tau\right]_f\right|^2}$$
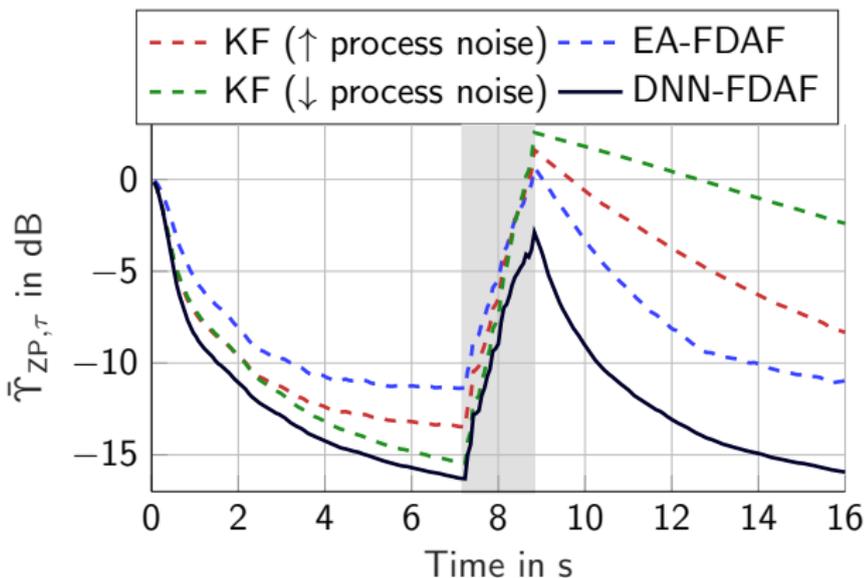
Static and frequency-independent
raw step-size

Haubner et al.: Deep Learning-Based Adaptation Control
Chair of Multimedia Communications and Signal Processing
May 2022
17 / 20

# System Identification Performance
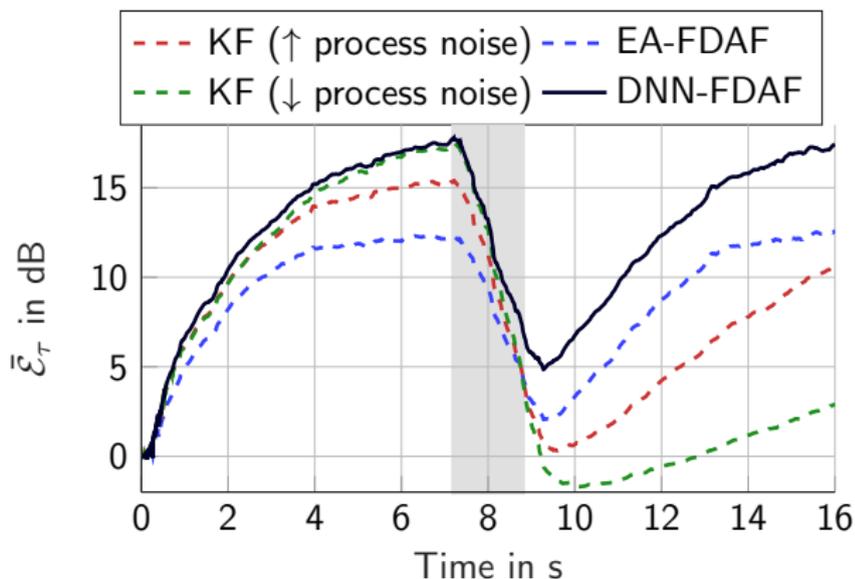


**Steady-state** and **tracking performance trade-off**

# System Identification Performance



**DNN-FDAF**: **Fast convergence** and **high steady-state** performance

# Echo Cancellation Performance



**DNN-FDAF**: Improved **echo cancellation** performance

# Conclusion

**Summary**

▶ Novel adaptation control for online system identification by using **DNN-based step-size inference**

▶ **End-to-end optimization** of DNN parameters w.r.t. average system identification performance

# Conclusion

**Summary**

► Novel adaptation control for online system identification by using **DNN-based step-size inference**

► **End-to-end optimization** of DNN parameters w.r.t. average system identification performance

**Outlook**

► **Joint control** of system identification and further parts of speech enhancement algorithms, e.g., spectral postfiltering

► Extension to **unsupervised system identification** applications

# Conclusion

**Summary**

▶ Novel adaptation control for online system identification by using **DNN-based step-size inference**

▶ **End-to-end optimization** of DNN parameters w.r.t. average system identification performance

**Outlook**

▶ **Joint control** of system identification and further parts of speech enhancement algorithms, e.g., spectral postfiltering

▶ Extension to **unsupervised system identification** applications

# Thank you for watching!