

Estimation of Bandwidth Extension Parameters in ITU-T G.729.1

— An “Add-on” to the Standard Decoder —

Bernd Geiser and Peter Vary

ETSI Workshop on Speech and Noise in Wideband Communication

Sophia Antipolis, France

May 22, 2007

Estimation of Bandwidth Extension Parameters in ITU-T G.729.1

— An “Add-on” to the Standard Decoder —

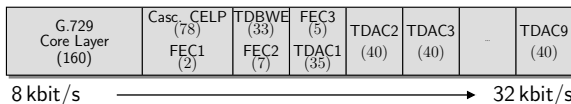
Outline

This work has been supported by Siemens AG, Munich, Germany.

- ITU-T G.729.1 — An Overview
- Proposed “Add-on”:
Estimation of TDBWE Parameters
- Results & Demonstration
- Discussion

ITU-T G.729.1 — An Overview

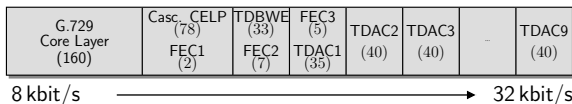
► Embedded Codec with Hierarchical Bitstream



- Core coder:** interoperable with G.729 (50 Hz – 4 kHz, 8 kbit/s)
- Cascade CELP:** [Massaloux et al. 2007]
improved narrowband quality (50 Hz – 4 kHz, 12 kbit/s)
- Bandwidth extension (TDBWE):** [Jax, Geiser et al. 2006]
good wideband speech quality (50 Hz – 7 kHz, 14 kbit/s)
- MDCT domain transform coding (TDAC):**
general audio capability (50 Hz – 7 kHz, up to 32 kbit/s)
- Additional information for *concealment of frame erasures* (FEC)
 - Bit rate adaptation is possible “on the fly”
 - G.729.1 is especially tailored for VoIP applications

ITU-T G.729.1 — An Overview

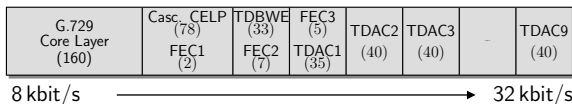
► Embedded Codec with Hierarchical Bitstream



- 1. Core coder:** interoperable with G.729 (50 Hz – 4 kHz, 8 kbit/s)
- 2. Cascade CELP:** [\[Massaloux et al. 2007\]](#)
improved narrowband quality (50 Hz – 4 kHz, 12 kbit/s)
- 3. Bandwidth extension (TDBWE):** [\[Jax, Geiser et al. 2006\]](#)
good wideband speech quality (50 Hz – 7 kHz, 14 kbit/s)
- 4. MDCT domain transform coding (TDAC):**
general audio capability (50 Hz – 7 kHz, up to 32 kbit/s)
- 5. Additional information for *concealment of frame erasures* (FEC)**
 - Bit rate adaptation is possible “on the fly”
 - G.729.1 is especially tailored for VoIP applications

ITU-T G.729.1 — An Overview

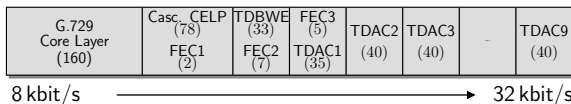
► Embedded Codec with Hierarchical Bitstream



- Core coder:** interoperable with G.729 (50 Hz – 4 kHz, 8 kbit/s)
- Cascade CELP:** [Massaloux *et al.* 2007]
improved narrowband quality (50 Hz – 4 kHz, 12 kbit/s)
- Bandwidth extension (TDBWE):** [Jax, Geiser *et al.* 2006]
good wideband speech quality (50 Hz – 7 kHz, 14 kbit/s)
- MDCT domain transform coding (TDAC):**
general audio capability (50 Hz – 7 kHz, up to 32 kbit/s)
- Additional information for *concealment of frame erasures* (FEC)
 - Bit rate adaptation is possible “on the fly”
 - G.729.1 is especially tailored for VoIP applications

ITU-T G.729.1 — An Overview

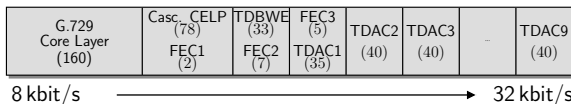
► Embedded Codec with Hierarchical Bitstream



- Core coder:** interoperable with G.729 (50 Hz – 4 kHz, 8 kbit/s)
- Cascade CELP:** [Massaloux *et al.* 2007]
improved narrowband quality (50 Hz – 4 kHz, 12 kbit/s)
- Bandwidth extension (TDBWE):** [Jax, Geiser *et al.* 2006]
good wideband speech quality (50 Hz – 7 kHz, 14 kbit/s)
- MDCT domain transform coding (TDAC):**
general audio capability (50 Hz – 7 kHz, up to 32 kbit/s)
- Additional information for *concealment of frame erasures* (FEC)
 - Bit rate adaptation is possible “on the fly”
 - G.729.1 is especially tailored for VoIP applications

ITU-T G.729.1 — An Overview

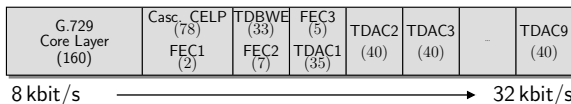
► Embedded Codec with Hierarchical Bitstream



- Core coder:** interoperable with G.729 (50 Hz – 4 kHz, 8 kbit/s)
- Cascade CELP:** [Massaloux *et al.* 2007]
improved narrowband quality (50 Hz – 4 kHz, 12 kbit/s)
- Bandwidth extension (TDBWE):** [Jax, Geiser *et al.* 2006]
good wideband speech quality (50 Hz – 7 kHz, 14 kbit/s)
- MDCT domain transform coding (TDAC):**
general audio capability (50 Hz – 7 kHz, up to 32 kbit/s)
- Additional information for *concealment of frame erasures* (FEC)
 - Bit rate adaptation is possible “on the fly”
 - G.729.1 is especially tailored for VoIP applications

ITU-T G.729.1 — An Overview

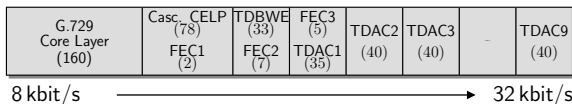
► Embedded Codec with Hierarchical Bitstream



- Core coder:** interoperable with G.729 (50 Hz – 4 kHz, 8 kbit/s)
- Cascade CELP:** [Massaloux *et al.* 2007]
improved narrowband quality (50 Hz – 4 kHz, 12 kbit/s)
- Bandwidth extension (TDBWE):** [Jax, Geiser *et al.* 2006]
good wideband speech quality (50 Hz – 7 kHz, 14 kbit/s)
- MDCT domain transform coding (TDAC):**
general audio capability (50 Hz – 7 kHz, up to 32 kbit/s)
- Additional information for *concealment of frame erasures* (FEC)
 - Bit rate adaptation is possible “on the fly”
 - G.729.1 is especially tailored for VoIP applications

ITU-T G.729.1 — An Overview

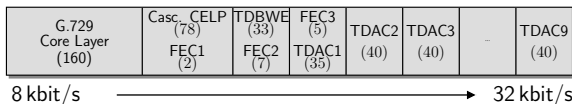
► Embedded Codec with Hierarchical Bitstream



- Core coder:** interoperable with G.729 (50 Hz – 4 kHz, 8 kbit/s)
- Cascade CELP:** [Massaloux *et al.* 2007]
improved narrowband quality (50 Hz – 4 kHz, 12 kbit/s)
- Bandwidth extension (TDBWE):** [Jax, Geiser *et al.* 2006]
good wideband speech quality (50 Hz – 7 kHz, 14 kbit/s)
- MDCT domain transform coding (TDAC):**
general audio capability (50 Hz – 7 kHz, up to 32 kbit/s)
- Additional information for *concealment of frame erasures* (FEC)
 - Bit rate adaptation is possible “on the fly”
 - G.729.1 is especially tailored for VoIP applications

ITU-T G.729.1 — An Overview

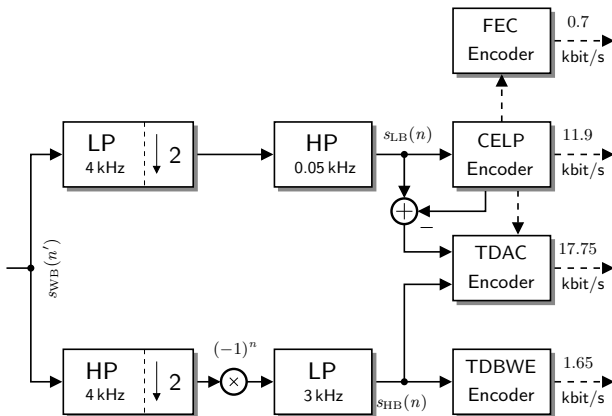
► Embedded Codec with Hierarchical Bitstream



- Core coder:** interoperable with G.729 (50 Hz – 4 kHz, 8 kbit/s)
- Cascade CELP:** [Massaloux *et al.* 2007]
improved narrowband quality (50 Hz – 4 kHz, 12 kbit/s)
- Bandwidth extension (TDBWE):** [Jax, Geiser *et al.* 2006]
good wideband speech quality (50 Hz – 7 kHz, 14 kbit/s)
- MDCT domain transform coding (TDAC):**
general audio capability (50 Hz – 7 kHz, up to 32 kbit/s)
- Additional information for *concealment of frame erasures* (FEC)
 - Bit rate adaptation is possible “on the fly”
 - G.729.1 is especially tailored for VoIP applications

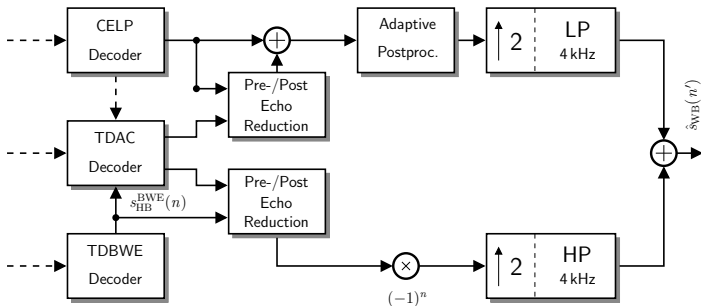
ITU-T G.729.1 — An Overview

► Encoder Signal Flow



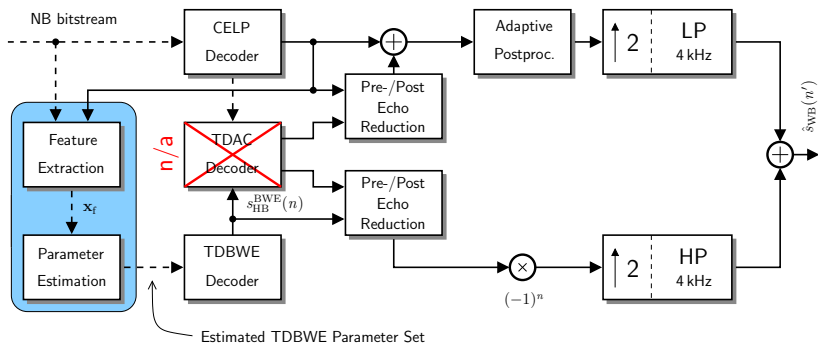
ITU-T G.729.1 — An Overview

► Decoder Signal Flow



New Extension: Estimation of TDBWE Parameters

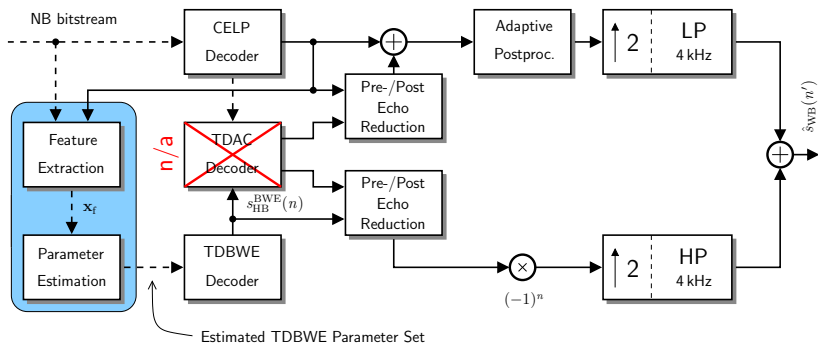
► Modified G.729.1 Decoder (8/12 kbit/s)



- MMSE estimation based on features from the low band bitstream
- Wideband speech at bit rates of 8 and 12 kbit/s

New Extension: Estimation of TDBWE Parameters

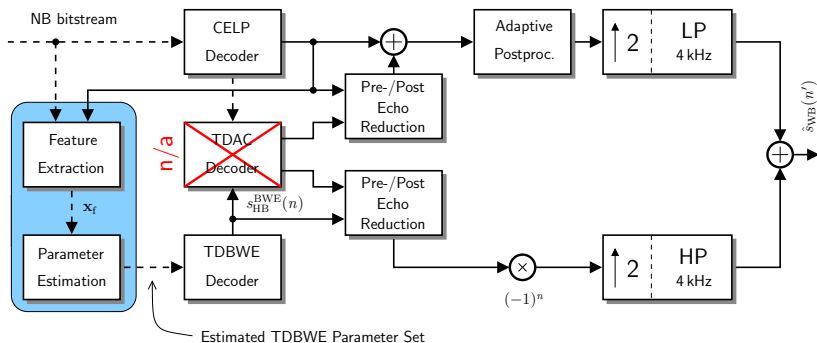
► Modified G.729.1 Decoder (8/12 kbit/s)



- MMSE estimation based on features from the low band bitstream
- Wideband speech at bit rates of 8 and 12 kbit/s

New Extension: Estimation of TDBWE Parameters

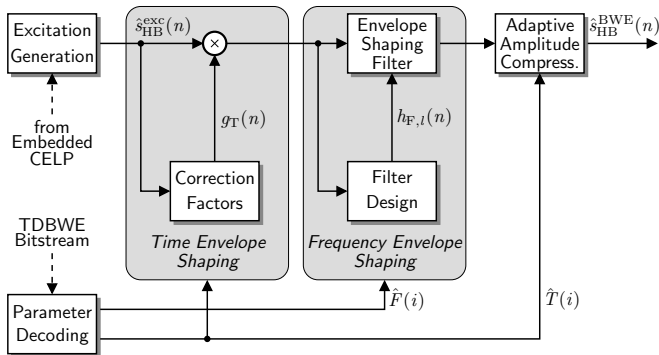
► Modified G.729.1 Decoder (8/12 kbit/s)



- MMSE estimation based on features from the low band bitstream
- Wideband speech at bit rates of 8 and 12 kbit/s

New Extension: Estimation of TDBWE Parameters

► Bandwidth Extension Algorithm (TDBWE)

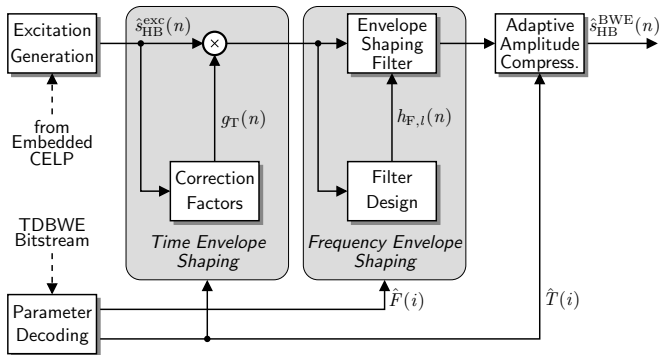


High band (4 kHz – 7 kHz) parameter set (per 20 ms frame):

- “Time envelope” T: 16 sub-frame energies (1.25 ms)
- “Frequency envelope” F: 12 sub-band energies (375 Hz)

New Extension: Estimation of TDBWE Parameters

► Bandwidth Extension Algorithm (TDBWE)

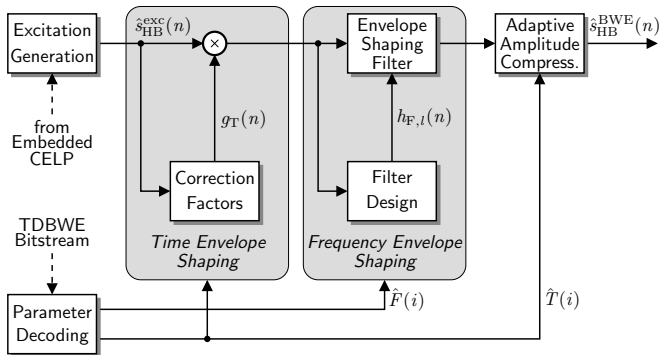


High band (4 kHz – 7 kHz) parameter set (per 20 ms frame):

- “Time envelope” T: 16 sub-frame energies (1.25 ms)
- “Frequency envelope” F: 12 sub-band energies (375 Hz)

New Extension: Estimation of TDBWE Parameters

► Bandwidth Extension Algorithm (TDBWE)

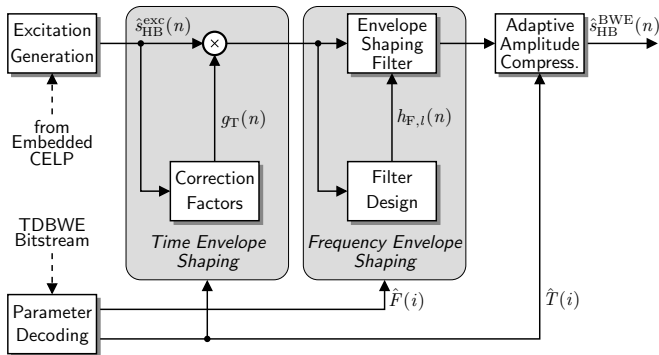


High band (4 kHz – 7 kHz) parameter set (per 20 ms frame):

- “Time envelope” **T**: 16 sub-frame energies (1.25 ms)
- “Frequency envelope” **F**: 12 sub-band energies (375 Hz)

New Extension: Estimation of TDBWE Parameters

▶ Bandwidth Extension Algorithm (TDBWE)



High band (4 kHz – 7 kHz) parameter set (per 20 ms frame):

- ▶ “Time envelope” **T**: 16 sub-frame energies (1.25 ms)
- ▶ “Frequency envelope” **F**: 12 sub-band energies (375 Hz)

New Extension: Estimation of TDBWE Parameters

▶ Narrowband Feature Vector \mathbf{x}_f per 20 ms Frame:

- ▶ Quantized spectral envelope (LSP vector) of the G.729 core codec
- ▶ Temporal envelope of the low band signal:
16 sub-frame energies (1.25 ms)

▶ Parameter Estimation

- ▶ MMSE estimation using *a priori* knowledge of 1st order [Jax, Vary 2003]

$$\tilde{\mathbf{y}}_{\text{MMSE}} = \sum_{\hat{\mathbf{y}}_i \in \mathcal{C}} \hat{\mathbf{y}}_i \cdot P(\hat{\mathbf{y}}_i | \mathbf{X}_f)$$

\mathbf{y} : "estimation quantity" (T or F)

$P(\hat{\mathbf{y}}_i | \mathbf{X}_f)$: "a posteriori" probabilities

\mathcal{C} : precomputed codebook

- ▶ Observation probability densities $p(\mathbf{x}_f | \hat{\mathbf{y}}_i)$ are approximated with GMMs (Gaussian Mixture Models)

New Extension: Estimation of TDBWE Parameters

▶ Narrowband Feature Vector \mathbf{x}_f per 20 ms Frame:

- ▶ Quantized spectral envelope (LSP vector) of the G.729 core codec
- ▶ Temporal envelope of the low band signal:
16 sub-frame energies (1.25 ms)

▶ Parameter Estimation

- ▶ MMSE estimation using *a priori* knowledge of 1st order [Jax, Vary 2003]

$$\tilde{\mathbf{y}}_{\text{MMSE}} = \sum_{\hat{\mathbf{y}}_i \in \mathcal{C}} \hat{\mathbf{y}}_i \cdot P(\hat{\mathbf{y}}_i | \mathbf{X}_f)$$

\mathbf{y} : "estimation quantity" (T or F)

$P(\hat{\mathbf{y}}_i | \mathbf{X}_f)$: "a posteriori" probabilities

\mathcal{C} : precomputed codebook

- ▶ Observation probability densities $p(\mathbf{x}_f | \hat{\mathbf{y}}_i)$ are approximated with GMMs (Gaussian Mixture Models)

New Extension: Estimation of TDBWE Parameters

▶ Narrowband Feature Vector \mathbf{x}_f per 20 ms Frame:

- ▶ Quantized spectral envelope (LSP vector) of the G.729 core codec
- ▶ Temporal envelope of the low band signal:
16 sub-frame energies (1.25 ms)

▶ Parameter Estimation

- ▶ MMSE estimation using *a priori* knowledge of 1st order [Jax, Vary 2003]

$$\tilde{\mathbf{y}}_{\text{MMSE}} = \sum_{\hat{\mathbf{y}}_i \in \mathcal{C}} \hat{\mathbf{y}}_i \cdot P(\hat{\mathbf{y}}_i | \mathbf{X}_f)$$

\mathbf{y} : “estimation quantity” (**T** or **F**)

$P(\hat{\mathbf{y}}_i | \mathbf{X}_f)$: “a posteriori” probabilities

\mathcal{C} : precomputed codebook

- ▶ Observation probability densities $p(\mathbf{x}_f | \hat{\mathbf{y}}_i)$ are approximated with GMMs (Gaussian Mixture Models)

New Extension: Estimation of TDBWE Parameters

▶ Narrowband Feature Vector \mathbf{x}_f per 20 ms Frame:

- ▶ Quantized spectral envelope (LSP vector) of the G.729 core codec
- ▶ Temporal envelope of the low band signal:
16 sub-frame energies (1.25 ms)

▶ Parameter Estimation

- ▶ MMSE estimation using *a priori* knowledge of 1st order [Jax, Vary 2003]

$$\tilde{\mathbf{y}}_{\text{MMSE}} = \sum_{\hat{\mathbf{y}}_i \in \mathcal{C}} \hat{\mathbf{y}}_i \cdot P(\hat{\mathbf{y}}_i | \mathbf{X}_f)$$

\mathbf{y} : “estimation quantity” (**T** or **F**)

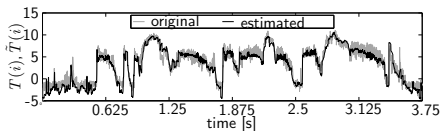
$P(\hat{\mathbf{y}}_i | \mathbf{X}_f)$: “a posteriori” probabilities

\mathcal{C} : precomputed codebook

- ▶ Observation probability densities $p(\mathbf{x}_f | \hat{\mathbf{y}}_i)$ are approximated with GMMs (Gaussian Mixture Models)

Results and Demonstration

▶ Example: Estimated vs. Original Time Envelope



▶ High Band Spectral Distortion (obtained with 8th order AR model)

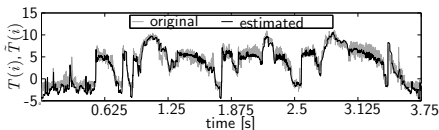
- ▶ **Proposed solution:** *estimated* TDBWE parameters, BWE *without* side information $\rightarrow \bar{d}_{\text{LSD}} \approx 5.05$ dB
- ▶ **Standardized solution:** *quantized* TDBWE parameters, BWE with 1.65 kbit/s side information $\rightarrow \bar{d}_{\text{LSD,quant.}} \approx 3.55$ dB

▶ Demonstration

| G.729.1 bit rate | 12 kbit/s | 12 kbit/s with BWE | 14 kbit/s |
|------------------|-----------|--------------------|-----------|
| Speaker 1 | 🔊 | 🔊 | 🔊 |
| Speaker 2 | 🔊 | 🔊 | 🔊 |

Results and Demonstration

▶ Example: Estimated vs. Original Time Envelope



▶ High Band Spectral Distortion (obtained with 8th order AR model)

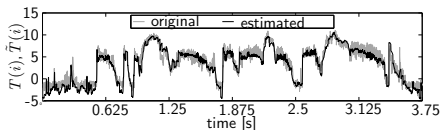
- ▶ **Proposed solution:** *estimated* TDBWE parameters, BWE *without* side information $\rightarrow \bar{d}_{\text{LSD}} \approx 5.05$ dB
- ▶ **Standardized solution:** *quantized* TDBWE parameters, BWE with 1.65 kbit/s side information $\rightarrow \bar{d}_{\text{LSD,quant.}} \approx 3.55$ dB

▶ Demonstration

| G.729.1 bit rate | 12 kbit/s | 12 kbit/s with BWE | 14 kbit/s |
|------------------|-----------|--------------------|-----------|
| Speaker 1 | 🔊 | 🔊 | 🔊 |
| Speaker 2 | 🔊 | 🔊 | 🔊 |

Results and Demonstration

▶ Example: Estimated vs. Original Time Envelope



▶ High Band Spectral Distortion (obtained with 8th order AR model)

- ▶ **Proposed solution:** *estimated* TDBWE parameters, BWE *without* side information $\rightarrow \bar{d}_{LSD} \approx 5.05$ dB
- ▶ **Standardized solution:** *quantized* TDBWE parameters, BWE with 1.65 kbit/s side information $\rightarrow \bar{d}_{LSD, quant.} \approx 3.55$ dB

▶ Demonstration

| G.729.1 bit rate | 12 kbit/s | 12 kbit/s with BWE | 14 kbit/s |
|------------------|-----------|--------------------|-----------|
| Speaker 1 | 🔊 | 🔊 | 🔊 |
| Speaker 2 | 🔊 | 🔊 | 🔊 |

Discussion

- ▶ Proposal: New *wideband* modes for G.729.1 at 8 and 12 kbit/s
- ▶ Estimation of TDBWE parameters based on the narrowband bitstream/signal → Encoder *not modified* (backwards compatible!)
- ▶ Band split frequency of 4 kHz leads to improved estimation performance (e.g. identification of fricative sounds)
- ▶ Satisfactory wideband speech quality is possible
- ▶ Efficient implementation through reuse of G.729.1 components

Application Scenarios:

- ▶ Increased speech quality and intelligibility at 8 and 12 kbit/s
- ▶ Heterogeneous conference scenario: constant acoustic bandwidth can be provided
- ▶ Treatment of bit rate switchings between 12 and 14 kbit/s

Discussion

- ▶ Proposal: New *wideband* modes for G.729.1 at 8 and 12 kbit/s
- ▶ Estimation of TDBWE parameters based on the narrowband bitstream/signal → Encoder *not modified* (backwards compatible!)
- ▶ Band split frequency of 4 kHz leads to improved estimation performance (e.g. identification of fricative sounds)
- ▶ Satisfactory wideband speech quality is possible
- ▶ Efficient implementation through reuse of G.729.1 components

Application Scenarios:

- ▶ Increased speech quality and intelligibility at 8 and 12 kbit/s
- ▶ Heterogeneous conference scenario: constant acoustic bandwidth can be provided
- ▶ Treatment of bit rate switchings between 12 and 14 kbit/s

Discussion

- ▶ Proposal: New *wideband* modes for G.729.1 at 8 and 12 kbit/s
- ▶ Estimation of TDBWE parameters based on the narrowband bitstream/signal → Encoder *not modified* (backwards compatible!)
- ▶ Band split frequency of 4 kHz leads to improved estimation performance (e.g. identification of fricative sounds)
- ▶ Satisfactory wideband speech quality is possible
- ▶ Efficient implementation through reuse of G.729.1 components

Application Scenarios:

- ▶ Increased speech quality and intelligibility at 8 and 12 kbit/s
- ▶ Heterogeneous conference scenario: constant acoustic bandwidth can be provided
- ▶ Treatment of bit rate switchings between 12 and 14 kbit/s

Discussion

- ▶ Proposal: New *wideband* modes for G.729.1 at 8 and 12 kbit/s
- ▶ Estimation of TDBWE parameters based on the narrowband bitstream/signal → Encoder *not modified* (backwards compatible!)
- ▶ Band split frequency of 4 kHz leads to improved estimation performance (e.g. identification of fricative sounds)
- ▶ Satisfactory wideband speech quality is possible
- ▶ Efficient implementation through reuse of G.729.1 components

Application Scenarios:

- ▶ Increased speech quality and intelligibility at 8 and 12 kbit/s
- ▶ Heterogeneous conference scenario: constant acoustic bandwidth can be provided
- ▶ Treatment of bit rate switchings between 12 and 14 kbit/s

Discussion

- ▶ Proposal: New *wideband* modes for G.729.1 at 8 and 12 kbit/s
- ▶ Estimation of TDBWE parameters based on the narrowband bitstream/signal → Encoder *not modified* (backwards compatible!)
- ▶ Band split frequency of 4 kHz leads to improved estimation performance (e.g. identification of fricative sounds)
- ▶ Satisfactory wideband speech quality is possible
- ▶ Efficient implementation through reuse of G.729.1 components

Application Scenarios:

- ▶ Increased speech quality and intelligibility at 8 and 12 kbit/s
- ▶ Heterogeneous conference scenario: constant acoustic bandwidth can be provided
- ▶ Treatment of bit rate switchings between 12 and 14 kbit/s

Discussion

- ▶ Proposal: New *wideband* modes for G.729.1 at 8 and 12 kbit/s
- ▶ Estimation of TDBWE parameters based on the narrowband bitstream/signal → Encoder *not modified* (backwards compatible!)
- ▶ Band split frequency of 4 kHz leads to improved estimation performance (e.g. identification of fricative sounds)
- ▶ Satisfactory wideband speech quality is possible
- ▶ Efficient implementation through reuse of G.729.1 components

Application Scenarios:

- ▶ Increased speech quality and intelligibility at 8 and 12 kbit/s
- ▶ Heterogeneous conference scenario: constant acoustic bandwidth can be provided
- ▶ Treatment of bit rate switchings between 12 and 14 kbit/s

Discussion

- ▶ Proposal: New *wideband* modes for G.729.1 at 8 and 12 kbit/s
- ▶ Estimation of TDBWE parameters based on the narrowband bitstream/signal → Encoder *not modified* (backwards compatible!)
- ▶ Band split frequency of 4 kHz leads to improved estimation performance (e.g. identification of fricative sounds)
- ▶ Satisfactory wideband speech quality is possible
- ▶ Efficient implementation through reuse of G.729.1 components

Application Scenarios:

- ▶ Increased speech quality and intelligibility at 8 and 12 kbit/s
- ▶ Heterogeneous conference scenario: constant acoustic bandwidth can be provided
- ▶ Treatment of bit rate switchings between 12 and 14 kbit/s

Discussion

- ▶ Proposal: New *wideband* modes for G.729.1 at 8 and 12 kbit/s
- ▶ Estimation of TDBWE parameters based on the narrowband bitstream/signal → Encoder *not modified* (backwards compatible!)
- ▶ Band split frequency of 4 kHz leads to improved estimation performance (e.g. identification of fricative sounds)
- ▶ Satisfactory wideband speech quality is possible
- ▶ Efficient implementation through reuse of G.729.1 components

Application Scenarios:

- ▶ Increased speech quality and intelligibility at 8 and 12 kbit/s
- ▶ Heterogeneous conference scenario: constant acoustic bandwidth can be provided
- ▶ Treatment of bit rate switchings between 12 and 14 kbit/s

References

ITU-T Rec. G.729.1, "G.729 based embedded variable bit-rate coder: An 8-32 kbit/s scalable wideband coder bitstream interoperable with G.729."

S. Ragot et al., "ITU-T G.729.1: An 8-32 kbit/s scalable coder interoperable with G.729 for wideband telephony and Voice over IP," in Proc. of IEEE ICASSP, Honolulu, Hawai'i, USA, Apr. 2007.

D. Massaloux et al., "An 8-12 kbit/s embedded CELP coder interoperable with ITU-T G.729 coder: First stage of the new G.729.1 standard", in Proc. of IEEE ICASSP, Honolulu, Hawai'i, USA, Apr. 2007.

P. Jax, B. Geiser, S. Schandl, H. Taddei, and P. Vary, "An embedded scalable wideband codec based on the GSM EFR codec," in Proc. of IEEE ICASSP, vol. 1, Toulouse, France, May 2006, pp. 5–8.

P. Jax and P. Vary, "On artificial bandwidth extension of telephone speech," Signal Processing, vol. 83, no. 8, pp. 1707–1719, Aug. 2003.

