# Recent Advances in Model-based Transform Audio Coding

Marie Oger, Stéphane Ragot

Flexcode Seminar Stockholm

October 17, 2007

- Flexible Coding for Heterogeneous Networks
- Objectives: develop **flexible source-channel coding algorithms**
  - More flexible than current, application-specific coders
  - Flexibility through online design, **generic source, channel and distortion models**
  - **Focus on audio**
- **http://www.flexcode.eu**



Nokia

Coordinator: KTH

Ericsson

RWTH Aachen

France Telecom

# Flexible <u>Source</u> Coding

- ## Approach
  - – Start with existing models: transform and linear-predictive coding
  - – "Flexcodize": analytic solutions $\rightarrow$ adaptive coding

- ## Tools for flexible coding include
  - – High-rate quantization theory
  - – Probability models (GMM, …) for quantizer design
    - Quantizer specification by equations
    - Estimate statistics for source
  - – Distortion measures using sensitivity matrix

- GMM (Gaussian Mixture Model) based LPC quantization [Subramaniam 01] [Samuelsson 01]
  - LPC coefficients or prediction error are modeled by a GMM
  - A mean-removed Karhunen-Loeve transform (KLT) and normalization by standard deviations is applied to LPC coefficients

- Autoregressive GMM for speech coding [Samuelsson 04]
  - Companded GMM for vector quantizers (CGMM-VQ)

- Generalized Gaussian model for image coding [Parisot 03]
  - Wavelet coding for image (EBWIC Coder)
  - Wavelet coefficients are modeled by a generalized Gaussian model

- ## Generalized Gaussian model
  - Definition
  - Example

- ## Proposed stack-run coding with model-based deadzone
  - Principle of stack-run coding
  - Rate control based on asymptotic bit allocation
  - Model-based optimization of deadzone
  - Objective & Subjective results
  - Delay & Complexity
  - Audio samples

- ## Latest developments: model-based bit plane coding
  - Principle
  - Preliminary results

- ## Conclusion & perspectives

# Generalized Gaussian pdf: definition

- The probability density function (pdf) of a zero-mean **generalized Gaussian variabl**e z of standard deviation σ is given by :

$$p_{\alpha,\sigma}(z) = \frac{A(\alpha)}{\sigma} e^{-|B(\alpha)z/\sigma|^{\alpha}}$$

where

$$A(\alpha) = \frac{\alpha B(\alpha)}{2\Gamma(1/\alpha)} \quad \text{and} \quad B(\alpha) = \sqrt{\frac{\Gamma(3/\alpha)}{\Gamma(1/\alpha)}}$$

with Γ(.) the Gamma function defined as

$$\Gamma(\alpha) = \int_{0}^{\infty} e^{-t} \, t^{\alpha+1} \, dt$$

- **The method used to estimate** α **is proposed by Mallat** [Mallat 89]

$$F(\alpha) = \frac{E(|z|)}{\sqrt{E(z^2)}} = \frac{\Gamma(2/\alpha)}{\sqrt{\Gamma(1/\alpha)\Gamma(3/\alpha)}}$$

So $\quad \hat{\alpha} = F^{-1}\left(\frac{\hat{m}_1}{\sqrt{\hat{m}_2}}\right) = F^{-1}\left(\frac{\sum_{i=1}^{n} z_i^2}{\sqrt{\sum_{i=1}^{n}|z_i|}}\right)$

signal segment (time)    spectrum (frequency)

Arithmetic-coded scalar quantization with model-based allocation

- Input/output signals sampled at 16 kHz
- Frame length of 20 ms with a lookahead of 25 ms (5 ms for LPC analysis and 20 ms for MDCT )
- Effective bandwith: 50-7000 Hz
- The perceptual weighting filter is defined as:

$$W(z) = \frac{A(z/\gamma)}{1-\beta z^{-1}} \quad \text{with} \quad \beta = 0.75 \quad \text{and} \quad \gamma = 0.92$$

- LPC coefficients quantized with a method based on GMM [Subramaniam 03] [Oger 06]
- MDCT implemented using the fast algorithm of [Duhamel 91]

Integer sequence → **Symbol mapping** → Quaternary sequence +,-,0,1 → **Contextual adaptive arithmetic coding** →

Context

- Stack-run coding is a **lossless coding method** representing integer sequences
  - Developed for wavelet image coding
- Adaptive arithmetic coding [Witten 87] using a **quaternary alphabet (0, 1, -, +)** and two contexts (one for "**runs**" and another for "**stacks**")
  - A run is a sequence of zeros
  - A stack is a non-zero signed integer

- Mapping rules for stack
  - The binary representation is unsigned
  - MSB is replaced by "+" if the coefficient is positive and "-" if it is negative.
  - The **absolute value is incremented by one**
  - The binary representation of "+4" is "+01" instead of "+00"
- The meanings of the symbol alphabet
  - "0" is used to signify a bit value of 0 in encoding of stack
  - "1" is used for bit value of 1 in stack, but it is not used for the MSB
  - "+" is used to represent the positive MSB of stack and for a bit value of **1 in representing run lengths**
  - "-" is used to represent the negative MSB of stack and for a bit value of **0 in representing run lengths**
- Mapping example for the sequence [0 0 0 +35 +4 0 0 0 0 0 0 0 0 0 0 -11]



| Run   | ++ |       |     | - + - |      |
|-------|----|-------|-----|-------|------|
| Stack |    | 00100+ | 10+ |       | 001- |

- Encoding of **N zero-mean independent variables $x_i$ of variances $\sigma_i^2$**

- In case of **high-resolution** the mean square error D [Gersho & Gray 93] is given by

$$D \approx \sum_{i=1}^{N} h_i \sigma_i^2 2^{-2b_i}$$

  where **$h_i$ is a function of the pdf of the variable $x_i$** and $b_i$ is the number of bits per sample used to code $x_i$

- For generalized Gaussian variables $x_i$ the factor $h_i$ is given by [Parisot 03] :

$$h_i = \frac{\Gamma\left(1/\alpha_i\right)^3}{3\alpha_i^2 \Gamma\left(3/\alpha_i\right)} e^{2/\alpha_i}$$

- Encoding of **N zero-mean variables $x_i$ of variances $\sigma_i^2$**
- The distortion D can be minimized by **Lagrangian techniques**:

$$J\left(b_i, \lambda\right) = D - \lambda\left(\sum_{i=1}^{N} b_i - B\right) \quad \Longrightarrow \quad \lambda_{opt} = 2\ln\left(2\right)\sum_{i=1}^{N} h_i \sigma_i^2 2^{-2b_i}$$

where B is the target bit rate

- Hence :

$$D_{opt} = \frac{\lambda_{opt}}{2\ln 2}$$

- In **case of high-resolution scalar uniform quantization** with step size q

$$D_{opt} = \frac{q_{opt}^2}{12} \quad \Longrightarrow \quad \boxed{q_{opt} = \sqrt{\frac{6\lambda_{opt}}{\ln 2}}}$$

# Model-based bit allocation: examples at 24 and 32 kbit/s

**Number of bits using high-rate estimated step size**



Bit allocation at 24 kbit/s (Budget = 430 bits/frame )

biais

Bit budget per frame



Bit allocation at 32 kbit/s (Budget = 590 bits/frame )

biais

- Biais due to mismatch high-rate assumption and use of context-based lossless coding instead of zero-entropy coding

➡ A bisection search is used in order to be within the bit budget constraint

- Encoding of **N zero-mean generalized Gaussian variables $x_i$ of variances $\sigma_i^2$**
- The distortion D is given by:

$$D(\alpha, z, q) = \frac{1}{\sigma^2} \int_{-z/2}^{z/2} x^2 p_{\sigma,\alpha}(x)\, dx + \frac{2}{\sigma^2} \sum_{m=1}^{+\infty} \int_{-z/2+(m-1)q}^{z/2+mq} (x - \hat{x}_m)^2 p_{\sigma,\alpha}(x)\, dx$$

- If the reconstruction level is set to centroid the distortion D is:

$$D(\alpha, z, q) = 1 - \sum_{m=1}^{+\infty} \frac{f_{1,m}\left(\alpha, \dfrac{z}{\sigma}, \dfrac{q}{\sigma}\right)^2}{f_{0,m}\left(\alpha, \dfrac{z}{\sigma}, \dfrac{q}{\sigma}\right)} \qquad \text{where} \quad f_{n,m}\left(\alpha, \frac{z}{\sigma}, \frac{q}{\sigma}\right) = \int_{z/2\sigma+(m-1)q/\sigma}^{z/2\sigma+mq/\sigma} x^n p_{1,\alpha}(x)\, dx$$

- If the reconstruction level is set to mid-value the distortion D is:

$$D(\alpha, z, q) = 1 + 2\sum_{m=1}^{+\infty} \left(\frac{1}{2}\frac{z}{\sigma} + \left(m - \frac{1}{2}\right)\frac{q}{\sigma}\right)^2 f_{0,m}\left(\alpha, \frac{z}{\sigma}, \frac{q}{\sigma}\right) - 4\sum_{m=1}^{+\infty} \left(\frac{1}{2}\frac{z}{\sigma} + \left(m - \frac{1}{2}\right)\frac{q}{\sigma}\right) f_{1,m}\left(\alpha, \frac{z}{\sigma}, \frac{q}{\sigma}\right)$$

- The bit rate R is given by:

$$R = -P(0)\log_2 P(0) - 2\sum_{m=1}^{+\infty} P(m)\log_2 P(m)$$

- With $P(m) = \int_{z/2+(m-1)q}^{z/2+mq} x^n p_\alpha(x)\,dx = f_{0,m}\left(\alpha, \frac{z}{\sigma}, \frac{q}{\sigma}\right)$

- So the bit rate R is given by:

$$R = -f_{0,0}\left(\alpha, \frac{z}{\sigma}\right)\log_2 f_{0,0}\left(\alpha, \frac{z}{\sigma}\right) - 2\sum_{m=1}^{+\infty} f_{0,m}\left(\alpha, \frac{z}{\sigma}, \frac{q}{\sigma}\right)\log_2 f_{0,m}\left(\alpha, \frac{z}{\sigma}, \frac{q}{\sigma}\right)$$

- Encoding of **N zero-mean generalized Gaussian variables $x_i$ of variances $\sigma_i^2$**
- The distortion D can be minimized by **Lagrangian techniques**:

$$J\left(z_i, q_i, \lambda\right) = \sum_{i=1}^{N} \sigma_i^2 D\left(\alpha_i, \frac{z_i}{\sigma_i}, \frac{q_i}{\sigma_i}\right) + \lambda\left(\sum_{i=1}^{N} a_i R\left(\alpha_i, \frac{z_i}{\sigma_i}, \frac{q_i}{\sigma_i}\right) - R_{\text{target}}\right)$$

$$\begin{cases} \dfrac{\dfrac{\partial D}{\partial \tilde{z}}\left(\alpha_i, \tilde{z}_i, \tilde{q}_i\right)}{\dfrac{\partial R}{\partial \tilde{z}}\left(\alpha_i, \tilde{z}_i, \tilde{q}_i\right)} = \dfrac{\dfrac{\partial D}{\partial \tilde{q}}\left(\alpha_i, \tilde{z}_i, \tilde{q}_i\right)}{\dfrac{\partial R}{\partial \tilde{q}}\left(\alpha_i, \tilde{z}_i, \tilde{q}_i\right)} \\[4ex] \dfrac{\dfrac{\partial D}{\partial \tilde{q}}\left(\alpha_i, \tilde{z}_i, \tilde{q}_i\right)}{\dfrac{\partial R}{\partial \tilde{q}}\left(\alpha_i, \tilde{z}_i, \tilde{q}_i\right)} = -\dfrac{\lambda a_i}{\sigma_i^2} \end{cases}$$

Optimization of the deadzone z

$$\sum_{i=1}^{N} a_i R\left(\alpha_i, \tilde{z}_i, \tilde{q}_i\right) = R_{\text{target}}$$

where $\alpha_i$, $z_i$, and $q_i$ are respectively the shape parameter, the deadzone and the stepsize

Scalar quantizer with reconstruction levels set to mid-value

Scalar quantizer with reconstruction levels set to optimal centroid (Lloyd-Max)

# Model-based deadzone optimization

(a) Female speech sample

(b) Shape parameter $\alpha$

(c) $\ln(q/\sigma)$ at 16 kbit/s (solid) and 32 kbit/s (dash)

(d) $z/q$ at 16 kbit/s (solid) and 32 kbit/s (dash)

- 16 clean music samples (4 types × 4 sentence-pairs) of 8 seconds
- 24 clean speech samples in French language (6 male and female speakers × 4 sentence-pairs) of 8 seconds
- Two AB test at 24 kbit/s : one for speech, another for music

A vs. B (A=G.722.1, B=proposed coder)

| | | | | |
|---|---|---|---|---|
| Music | 3% | 21% | 23% | 39% | 14% |
| Speech | 3% | 27% | 22% | 30% | 18% |

Legend:
- ▣ A is better than B
- ▨ A is slightly better than B
- ■ no preference
- ▨ B is slightly better than A
- ▨ B is better than A

0%    25%    50%    75%    100%

- 8 expert listeners
- The stack-run coding (z=q) was preferred:
  - In 53% cases for music
  - In 48% cases for speech
- Informal listening tests at 32 kbit/s
- **Stack-run coding (z=q) is better than  ITU-T G.722.1 at 24 kbit/s and equivalent at 32 kbit/s**

- 24 clean speech samples in French language (6 male and female speakers × 4 sentence-pairs) of 8 s
- Proposed coder: predictive MDCT coder with stack-run coding
- Results are presented with noise injection (injection similar to 3GPP AMR-WB+)
- These **objective results suggest that the inclusion of a dead-zone improves the performance**

- 20 clean speech samples in French language (5 male and female speakers $\times$ 4 sentence-pairs) of 8 seconds

A vs. B (A=z=z$_{opt}$, B=z=q)

| | | | | |
|---|---|---|---|---|
| Speech | 2% | 48% | 27% | 23% 0% |

0%  25%  50%  75%  100%

- A is better than B
- A is sligtly better than B
- no preference
- B is sligtly better than A
- B is better than A

- One AB test at 24 kbit/s for speech:
  - 9 expert listeners
  - Stack-run coding with $z=z_{opt}$ was preferred in 50% cases for speech
- Informal listening tests at 32 kbit/s
- **Stack-run coding with $z=z_{opt}$ is better than stack-run coding with z=q at low bitrate and is equivalent at high bitrate.**

- Algorithmic delay:
  - 45 ms (20 ms for the frame, 20ms for the MDCT and 5 ms for the lookahead) for the stack-run coder
  - 40 ms for the ITU-T G.722.1
- The computational complexity of ITU-T G.722.1 is very low
- The computational complexity of the stack-run coder is higher due to the use of bisection search for bit rate matching
  - Stack-run coding is performed several times per frame
- Storage requirements for the stack-run coder are low
  - Parameters for GMM based LPC quantization
  - MDCT tables (can be computed on lines)

Comparison between stack-run coding and ITU-T G.722.1

|  | Stack-run coding with z=q | Stack-run coding with z=$z_{opt}$ | ITU-T G.722.1 |
|---|---|---|---|
| Music at 24 kbit/s | 🔊 | 🔊 | 🔊 |
| Music at 32 kbit/s | 🔊 | 🔊 | 🔊 |
| Speech at 24 kbit/s | 🔊 | 🔊 | 🔊 |
| Speech at 32 kbit/s | 🔊 | 🔊 | 🔊 |

Bit-plane-coded scalar quantization with model-based allocation and probabilities



- Principle: replace stack-run coding by bit plane coding.
- Implicit rate control

➡ Computational complexity is much more lower than stack-run coding

Normalized MDCT spectrum



| Y= | -2 | +5 | -3 | ……... | +1 | 0 | +1 | |
|---|---|---|---|---|---|---|---|---|
| | 1 | 0 | 1 | ……... | 0 | 0 | 0 | → Bit plane of signs |
| $2^2$ | 0 | 1 | 0 | 0……...0 | 0 | 0 | 0 | → $P_2$=MSB |
| $2^1$ | 1 | 0 | 1 | 0……..0 | 0 | 0 | 0 | → |
| $2^0$ | 0 | 1 | 1 | 0……..1 | 1 | 0 | 1 | → $P_0$=LSB |

K=3

- The normalized MDCT spectrum $X_{pre}(k)$ is scalar quantized and we get an integer sequence $Y(k)$.

- This integer sequence is decomposed in binary format.

- The symbol probabilities in bit planes are estimated on the model of the pdf of $X_{pre}(k)$

# Objective quality results

- 24 clean speech samples in French language (6 male and female speakers × 4 sentence-pairs) of 8 s
- Proposed coder: predictive MDCT coder with bit-plane coding
- Results are presented without noise injection
- These **objective results suggest that the speech quality of the proposed coder with model-based initialization of symbol probabilities is equivalent to reference coders at high bitrate**

- We proposed a predictive MDCT coder with generalized Gaussian modeling for wideband speech and audio signals

- Generalized Gaussian modeling is used to:
  - Estimate the optimal step size
  - Optimize the deadzone
  - Estimate symbol probabitilies in bit planes

- Next step: Include sensitivity matrix into model-based coder
  - Linear-predictive filter $\rightarrow$ signal-adaptive transform

*FlexCode*

# Thank you!