

UNIFIED PRINCIPLE OF $\Sigma\Delta$ MODULATION AS QUANTIZATION TECHNIQUE FOR OVERCOMPLETE EXPANSIONS

Thao Nguyen

Dept. of Electrical Engineering

City University of New York

DISCRETIZATION OF SIGNAL

$\{\boldsymbol{\Phi}_n\}_{n \in \mathbb{Z}}$: generating family of vectors spanning input space

$$\mathbf{x} = \sum_{n \in \mathbb{Z}} x_n \cdot \boldsymbol{\Phi}_n$$

$$\hat{\mathbf{x}} = \sum_{n \in \mathbb{Z}} q_n \cdot \boldsymbol{\Phi}_n \quad \text{where } q_n \in \underbrace{\{l_1, l_2, \dots, l_N\}}_{\text{quantization levels}}$$

$$\text{error } \mathbf{x} - \hat{\mathbf{x}} = \sum_{n \in \mathbb{Z}} (x_n - q_n) \cdot \boldsymbol{\Phi}_n$$

USE OF REDUNDANCY IN $\Sigma\Delta$ MODULATION

Form vectors $\mathbf{r}_k := \boldsymbol{\varphi}_k - \sum_{n \neq k} c_{n,k} \boldsymbol{\varphi}_n, \quad k \in \mathbb{Z}$

General notation $\mathbf{r}_k := \sum_{n \in \mathbb{Z}} d_{n,k} \boldsymbol{\varphi}_n, \quad k \in \mathbb{Z}, \quad \text{with } d_{k,k} = 1$

$D := \{d_{n,k}\}_{n,k \in \mathbb{Z}}$: redundancy operator

$$\mathbf{x} - \hat{\mathbf{x}} = \sum_{n \in \mathbb{Z}} (x_n - q_n) \cdot \boldsymbol{\varphi}_n \quad \Leftrightarrow \quad \mathbf{x} - \hat{\mathbf{x}} = \sum_{k \in \mathbb{Z}} u_k \cdot \mathbf{r}_k$$

where $x_n - q_n = \sum_{k \in \mathbb{Z}} d_{n,k} \cdot u_k$

PROOF

$$\mathbf{x} - \hat{\mathbf{x}} = \sum_{n \in \mathbb{Z}} (x_n - q_n) \cdot \boldsymbol{\varphi}_n$$

$$x_n - q_n = \sum_{k \in \mathbb{Z}} d_{n,k} \cdot u_k$$

$$\mathbf{r}_k = \sum_{n \in \mathbb{Z}} d_{n,k} \cdot \boldsymbol{\varphi}_n$$

$$\mathbf{x} - \hat{\mathbf{x}} = \sum_{n \in \mathbb{Z}} \underbrace{\left(\sum_{k \in \mathbb{Z}} d_{n,k} \cdot u_k \right)}_{x_n - q_n} \cdot \boldsymbol{\varphi}_n = \sum_{k \in \mathbb{Z}} u_k \cdot \underbrace{\left(\sum_{n \in \mathbb{Z}} d_{n,k} \cdot \boldsymbol{\varphi}_n \right)}_{\mathbf{r}_k}$$

$$\mathbf{x} - \hat{\mathbf{x}} = \sum_{k \in \mathbb{Z}} u_k \cdot \mathbf{r}_k$$

DIFFERENTIATION EXAMPLE

Form vectors

$$\mathbf{r}_k := \boldsymbol{\varphi}_k - \boldsymbol{\varphi}_{k+1}$$

$$\mathbf{x} - \hat{\mathbf{x}} = \sum_{n \in \mathbb{Z}} (x_n - q_n) \cdot \boldsymbol{\varphi}_n \quad \Leftrightarrow \quad \mathbf{x} - \hat{\mathbf{x}} = \sum_{k \in \mathbb{Z}} u_k \cdot \mathbf{r}_k$$

where $x_n - q_n = u_n - u_{n-1}$

DIFFERENTIATION EXAMPLE

Form vectors

$$\mathbf{r}_k := \boldsymbol{\varphi}_k - \boldsymbol{\varphi}_{k+1}$$

$$\mathbf{r}_k := \sum_{n \in \mathbb{Z}} d_{n,k} \boldsymbol{\varphi}_n, \quad k \in \mathbb{Z},$$

where $d_{n,k} = d_{n-k}$ with $d_n := \delta_n - \delta_{n+1}$

$$D := \{d_{n-k}\}_{n,k \in \mathbb{Z}}$$

particular case of “convolutional” or “shift-invariant” operator

PRINCIPLES OF $\Sigma\Delta$ MODULATION

Choose redundancy operator $D = \{d_{n,k}\}_{n,k \in \mathbb{Z}}$ invertible and

such that $\mathbf{r}_k := \sum_{n \in \mathbb{Z}} d_{n,k} \boldsymbol{\varphi}_n$ are “small”

Find quantized sequence $q_n \in \{l_1, l_2, \dots, l_N\}$ so that equation

$$x_n - q_n = \sum_{k \in \mathbb{Z}} d_{n,k} \cdot u_k$$

yields bounded and “small” solution in u_k

Hopefully, $\mathbf{x} - \hat{\mathbf{x}} = \sum_{n \in \mathbb{Z}} u_n \cdot \mathbf{r}_n$ will be “small”

PRINCIPLES OF $\Sigma\Delta$ MODULATION

Choose redundancy operator $D = \{d_{n,k}\}_{n,k \in \mathbb{Z}}$ invertible and

such that $\mathbf{r}_k := \sum_{n \in \mathbb{Z}} d_{n,k} \boldsymbol{\varphi}_n$ are “small”

Hopefully, $\mathbf{x} - \hat{\mathbf{x}} = \sum_{n \in \mathbb{Z}} u_n \cdot \mathbf{r}_n$ will be “small”

SPACE NORM

Main case of study: shift-invariant space of functions

$\boldsymbol{\varphi}_n = \varphi(t - n\tau)$ where $\varphi(t)$ and τ are given

• $\left(\begin{array}{l} \|\mathbf{x}\|_2^2 = \int_{-\infty}^{\infty} |x(t)|^2 dt \quad \text{however} \quad \hat{\mathbf{x}} = \sum_{n \in \mathbb{Z}} q_n \cdot \boldsymbol{\varphi}_n \notin L^2 \\ \text{use of L2 norm requires special treatment, see [O.Yilmaz, 2004]} \end{array} \right)$

• $\|\mathbf{x}\|_{\infty} = \sup_{t \in \mathbb{R}} |x(t)|$

• $\text{MSE}(\mathbf{x}) = \lim_{T \rightarrow \infty} \frac{1}{T} \int_{-T}^T |x(t)|^2 dt$

SPACE NORM

Main case of study: shift-invariant space of functions

$\varphi_n = \varphi(t - n\tau)$ where $\varphi(t)$ and τ are given

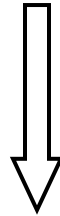
•

•
$$\|\mathbf{x}\|_{\infty} = \sup_{t \in \mathbf{R}} |x(t)|$$

•
$$\text{MSE}(\mathbf{x}) = \lim_{T \rightarrow \infty} \frac{1}{T} \int_{-T}^T |x(t)|^2 dt$$

SUP NORM BOUND

$$\mathbf{x} - \hat{\mathbf{x}} = \sum_{n \in \mathbb{Z}} u_n \cdot \mathbf{r}_n$$



$$|x(t) - \hat{x}(t)| \leq \|u\|_{\infty} \cdot \sum_{k \in \mathbb{Z}} |r_k(t)|$$

[I.Daubechies & R.DeVore, 2003]

DIFFERENTIATION

$$D = \{d_{n-k}\}_{n,k \in \mathbb{Z}} \quad \text{with} \quad d_n = \delta_n - \delta_{n-1}$$

$$\mathbf{r}_k := \boldsymbol{\varphi}_k - \boldsymbol{\varphi}_{k+1} \quad \text{with} \quad \boldsymbol{\varphi}_n = \varphi(t - n\tau)$$



$$\sum_{k \in \mathbb{Z}} |r_k(t)| \leq \|\varphi'\|_{L^1}$$

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_{\infty} \leq \|u\|_{\infty} \cdot \|\varphi'\|_{L^1}$$

[I.Daubechies & R.DeVore, 2003]

APPLICATION TO BANDLIMITED INPUTS

Assumption : $x(t)$ lowpass of bandwidth $\Omega_0 = \frac{2\pi}{\tau_0}$

Then $x(t) := \sum_{n \in \mathbb{Z}} x_n \cdot \varphi(t - n\tau)$

where $x_n = x(n\tau)$ and $\varphi(t) = \tau \operatorname{sinc}_{\tau_0}(t)$

with $\operatorname{sinc}_{\tau_0}(t) := \frac{\sin(\pi t / \tau_0)}{\pi t}$ and $\tau < \tau_0 := \frac{2\pi}{\Omega_0} \rightarrow R = \frac{\tau_0}{\tau} = \text{redundancy}$

↑
up to some relaxation in
the frequency domain, see *

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_{\infty} \leq \tau \cdot \|u\|_{\infty} \cdot \left\| \operatorname{sinc}_{\tau_0} \right\|_{L^1}$$

[* I.Daubechies & R.DeVore, 2003]

m^{th} ORDER DIFFERENTIATION

$$D = \{d_{n-k}\}_{n,k \in \mathbb{Z}} \quad \text{with} \quad d_n = (\delta_n - \delta_{n-1})^{(m)}$$

$$\mathbf{r}_k := \sum_{n \in \mathbb{Z}} d_{n-k} \boldsymbol{\varphi}_n \quad \text{with} \quad \boldsymbol{\varphi}_n = \varphi(t - n\tau)$$



$$\sum_{k \in \mathbb{Z}} |r_k(t)| \leq \tau^{m-1} \|\varphi^{(m)}\|_{L^1}$$

$$|x(t) - \hat{x}(t)| \leq \|u\|_{\infty} \cdot \sum_{k \in \mathbb{Z}} |r_k(t)|$$

[I.Daubechies & R.DeVore, 2003]

APPLICATION TO BANDLIMITED INPUTS

$$D = \{d_{n-k}\}_{n,k \in \mathbb{Z}} \quad \text{with} \quad d_n = (\delta_n - \delta_{n-1})^{(m)}$$

$$\sum_{k \in \mathbb{Z}} |r_k(t)| \leq \tau^{m-1} \left\| \text{sinc}_{\tau_0}^{(m)} \right\|_{L^1}$$

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_{\infty} \leq \tau^m \cdot \|u\|_{\infty} \cdot \left\| \text{sinc}_{\tau_0}^{(m)} \right\|_{L^1}$$

[I.Daubechies & R.DeVore, 2003]

SPACE NORM

Main case of study: shift-invariant space of functions

$\varphi_n = \varphi(t - n\tau)$ where $\varphi(t)$ and τ are given

•

•
$$\|\mathbf{x}\|_{\infty} = \sup_{t \in \mathbf{R}} |x(t)|$$

•
$$\text{MSE}(\mathbf{x}) = \lim_{T \rightarrow \infty} \frac{1}{T} \int_{-T}^T |x(t)|^2 dt$$

MSE ANALYSIS

Main case of study: shift-invariant space of functions
shift-invariant operator D

$\boldsymbol{\varphi}_n = \varphi(t - n\tau)$ where $\varphi(t)$ and τ are given

$\Rightarrow \mathbf{r}_k = r(t - k\tau)$ where $r(t) = \sum_{n \in \mathbb{Z}} d_n \varphi(t - n\tau)$

Theorem:

$$\mathbf{x} - \hat{\mathbf{x}} = \sum_{k \in \mathbb{Z}} u_k \cdot \mathbf{r}_k \quad \Rightarrow \quad \text{MSE}(\mathbf{x} - \hat{\mathbf{x}}) = \sum_{k \in \mathbb{Z}} a_k s_k$$

where

$$a_k := \lim_{N \rightarrow \infty} \frac{1}{2N} \sum_{n=-N}^N u_n u_{n+k} \quad \text{and} \quad s_k := \frac{1}{\tau} \int_{\mathbb{R}} r(t) r(t - k\tau) dt$$

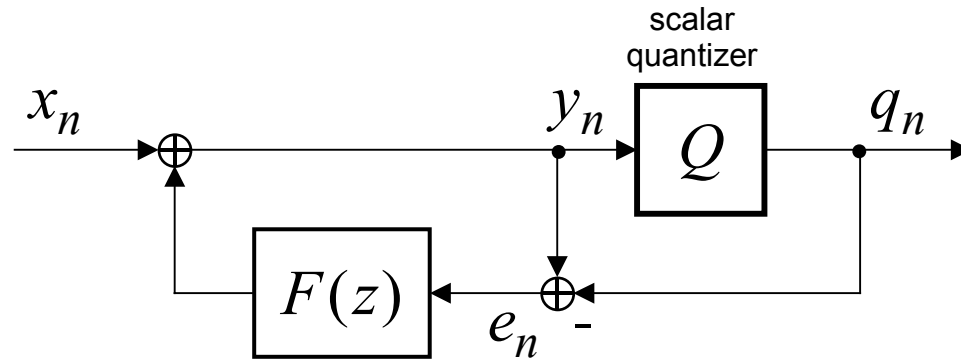
STATISTICAL MSE ANALYSIS

$$E\{\text{MSE}(\mathbf{x} - \hat{\mathbf{x}})\} = \sum_{k \in \mathbb{Z}} \bar{a}_k s_k$$

where

$$\bar{a}_k := E\{u_n u_{n+k}\} \quad \text{and} \quad s_k := \frac{1}{\tau} \int_{\mathbf{R}} r(t) r(t - k\tau) dt$$

“WHITE NOISE” BEHAVIOR



$$u_n = g_n * e_n$$

$$E\{e_n e_{n+k}\} \xrightarrow{\text{large \#bits}} \sigma_e^2 \delta_k$$

$$\Rightarrow \bar{a}_k := E\{u_n u_{n+k}\} \xrightarrow{\text{large \#bits}} \bar{a}_0 \delta_k$$

“WHITE NOISE” MODEL

$$E\{\text{MSE}(\mathbf{x} - \hat{\mathbf{x}})\} = \frac{\bar{a}_0}{\tau} \|\mathbf{r}\|_2^2$$

$$\bar{a}_k := E\{u_n u_{n+k}\} = \bar{a}_0 \delta_k$$

$$s_0 = \frac{1}{\tau} \int_{\mathbf{R}} |r(t)|^2 dt = \frac{1}{\tau} \|\mathbf{r}\|_2^2$$

FIR DESIGN OF D UNDER “WHITE NOISE” MODEL

Assume constraint

$$D = \{d_{n-k}\}_{n,k \in \mathbb{Z}} \text{ with } \{d_n\}_{n \in \mathbb{Z}} = \underbrace{\{d_0, d_1, \dots, d_m\}}_{=1}$$

$$r(t) = \sum_{n \in \mathbb{Z}} d_n \varphi(t - n\tau) \quad \Rightarrow \quad \mathbf{r} = \boldsymbol{\varphi}_0 + d_1 \boldsymbol{\varphi}_1 + \dots + d_m \boldsymbol{\varphi}_m$$

$$E\{\text{MSE}(\mathbf{x} - \hat{\mathbf{x}})\} = \frac{\bar{a}_0}{\tau} \|\mathbf{r}\|_2^2$$

is minimized when d_1, \dots, d_m are chosen such that

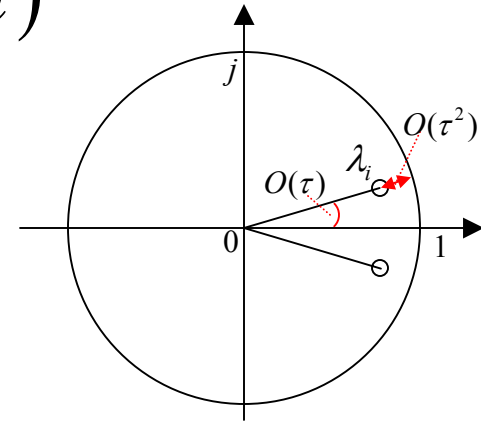
$$d_1 \boldsymbol{\varphi}_1 + \dots + d_m \boldsymbol{\varphi}_m = -\text{proj}_{\langle \boldsymbol{\varphi}_1, \dots, \boldsymbol{\varphi}_m \rangle}(\boldsymbol{\varphi}_0)$$

BANDLIMITED CASE

$$\varphi(t) = \tau \operatorname{sinc}_{\tau_0}(t)$$

White noise model optimization

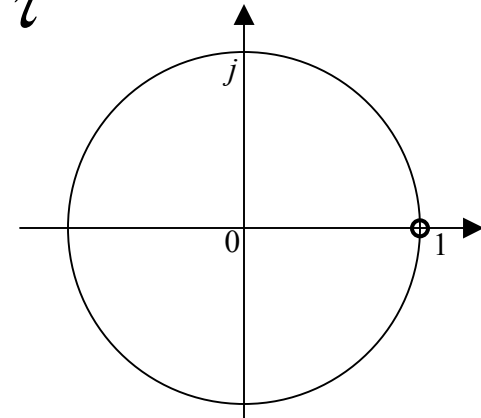
$$\Rightarrow D(z) = \prod_{i=1}^m (1 - \lambda_i z^{-1})$$



$$E\{\text{MSE}(\mathbf{x} - \hat{\mathbf{x}})\} \approx \alpha \tau^{2m+1}$$

Differentiation operator

$$D(z) = (1 - z^{-1})^m$$



$$E\{\text{MSE}(\mathbf{x} - \hat{\mathbf{x}})\} \approx \beta \tau^{2m+1}$$

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_{\infty}^2 \leq \gamma \tau^{2m}$$

PRINCIPLES OF $\Sigma\Delta$ MODULATION

Choose redundancy operator $D = \{d_{n,k}\}_{n,k \in \mathbb{Z}}$ invertible and

such that $\mathbf{r}_k := \sum_{n \in \mathbb{Z}} d_{n,k} \boldsymbol{\varphi}_n$ are “small”

Find quantized sequence $q_n \in \{l_1, l_2, \dots, l_N\}$ so that equation

$$x_n - q_n = \sum_{k \in \mathbb{Z}} d_{n,k} \cdot u_k$$

yields bounded and “small” solution in u_k

Hopefully, $\mathbf{x} - \hat{\mathbf{x}} = \sum_{n \in \mathbb{Z}} u_n \cdot \mathbf{r}_n$ will be “small”

PRINCIPLES OF $\Sigma\Delta$ MODULATION

Find quantized sequence $q_n \in \{l_1, l_2, \dots, l_N\}$ so that equation

$$x_n - q_n = \sum_{k \in \mathbb{Z}} d_{n,k} \cdot u_k$$

yields bounded and “small” solution in u_k

CAUSAL AND TIME-INVARIANT FRAMEWORK

- $D = \{d_{n,k}\}_{n,k \in \mathbb{Z}}$ such that $d_{n,k} = d_{n-k}$ with $d_n = 0, \forall n < 0$ and $d_0 = 1$
- q_n must be decided at instant n

Find quantized sequence $q_n \in \{l_1, l_2, \dots, l_N\}$ so that equation

$$x_n - q_n = d_n * u_n$$

yields bounded and “small” solution in u_k

CAUSAL AND TIME-INVARIANT FRAMEWORK

- $D = \{d_{n,k}\}_{n,k \in \mathbb{Z}}$ such that $d_{n,k} = d_{n-k}$ with $d_n = 0, \forall n < 0$ and $d_0 = 1$
- q_n must be decided at instant n

Find quantized sequence $q_n \in \{l_1, l_2, \dots, l_N\}$ so that equation

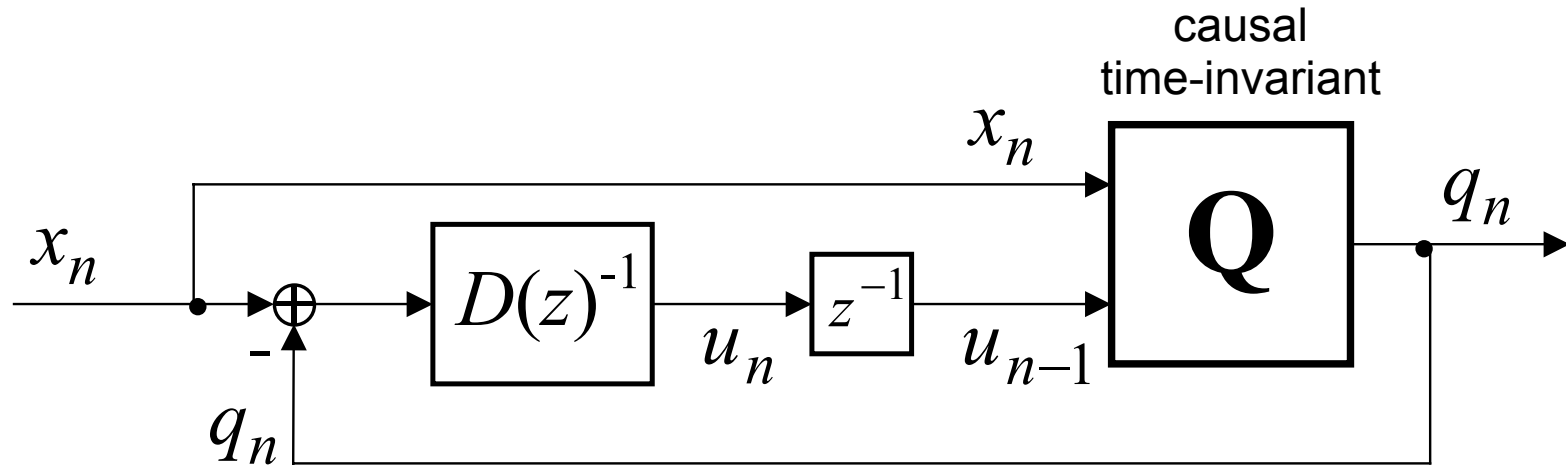
$$u_n = -\sum_{k>0} d_k \cdot u_{n-k} + (x_n - q_n)$$

yields bounded and “small” solution in u_k

$$q_n = \mathbf{Q}(x_n, x_{n-1}, \dots; u_{n-1}, u_{n-2}, \dots) \\ \in \{l_1, l_2, \dots, l_N\}$$

CAUSAL AND TIME-INVARIANT RESOLUTION

dynamical system

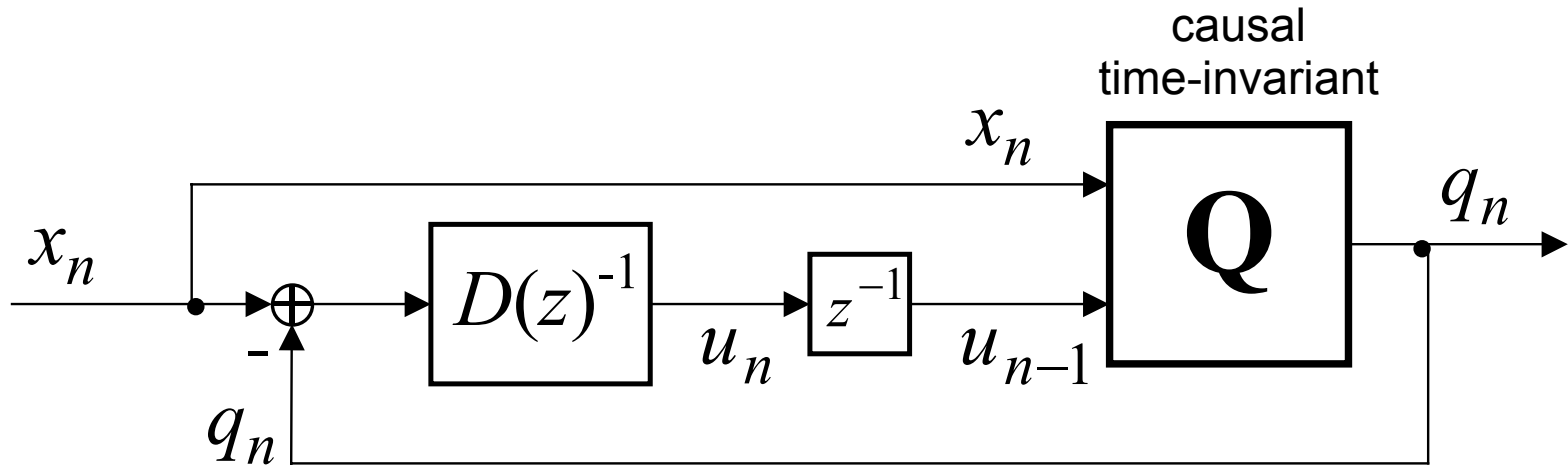


$$u_n = - \sum_{k>0} d_k \cdot u_{n-k} + (x_n - q_n)$$

$$q_n = \mathbf{Q}(x_n, x_{n-1}, \dots; u_{n-1}, u_{n-2}, \dots)$$

STABLE ONE-BIT SCHEME

[Daubechies & DeVore (2003)]

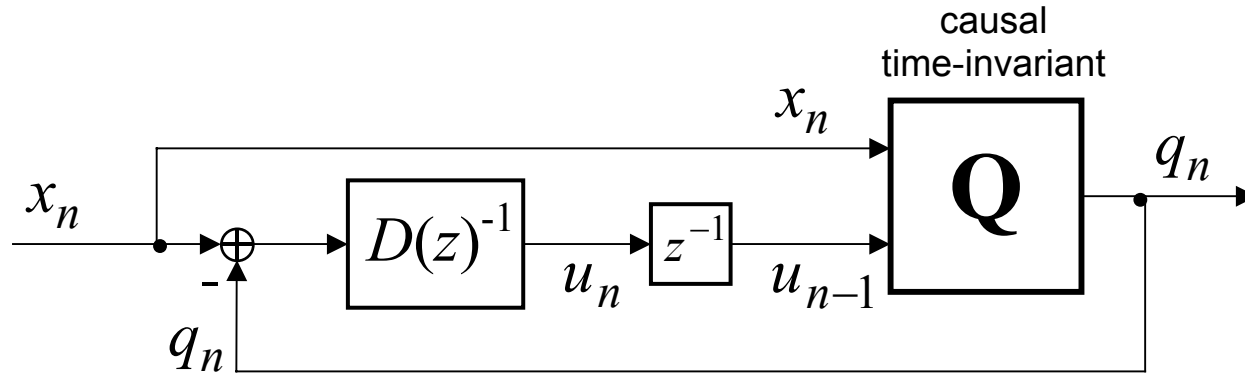


$$D = \{d_{n-k}\}_{n,k \in \mathbb{Z}} \text{ such that } d_n = (\delta_n - \delta_{n-1})^{(m)} \quad (D(z) = (1 - z^{-1})^m)$$

$$|x_n| \leq \frac{\Delta}{2} - \varepsilon$$

$$q_n = \mathbf{Q}(x_n; u_{n-1}, u_{n-2}, \dots, u_{n-m}) \in \left\{ +\frac{\Delta}{2}, -\frac{\Delta}{2} \right\}$$

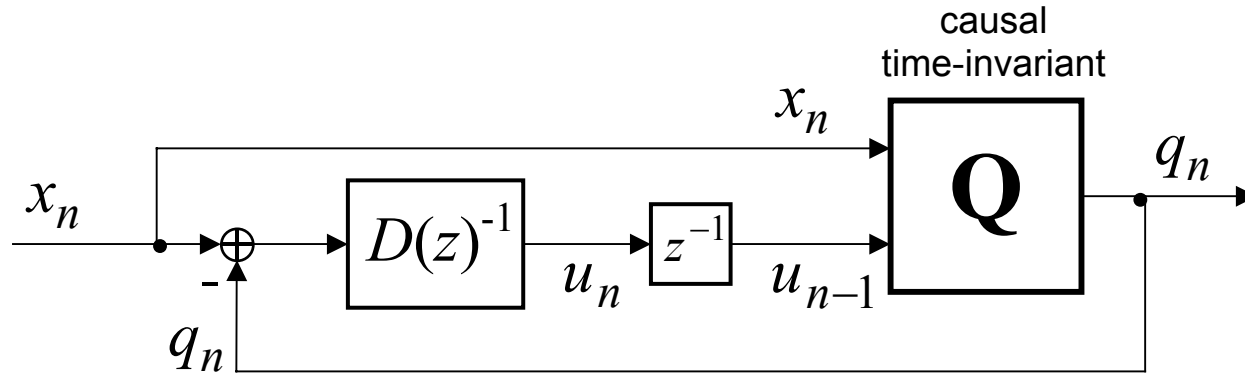
RESTRICTED SCHEME



restriction: $\mathbf{Q}(x_n, x_{n-1}, \dots; u_{n-1}, u_{n-2}, \dots) = Q(x_n + c_n * u_n)$

where c_n is strictly causal and Q is a **scalar** quantizer

RESTRICTED SCHEME



restriction:
$$\mathbf{Q}(x_n, x_{n-1}, \dots; u_{n-1}, u_{n-2}, \dots) = \mathbf{Q}(\underbrace{x_n + c_n * u_n}_{y_n})$$

$$q_n = Q(y_n)$$

$$x_n - q_n = d_n * u_n$$

$$y_n = x_n + c_n * u_n$$

RESTRICTED SCHEME

Define

$$e_n := y_n - q_n$$

Perform change of state variable

$$e_n = (d_n + c_n) * u_n$$

$$q_n = Q(y_n)$$

$$y_n - q_n = (d_n + c_n) * u_n$$

$$y_n = x_n + c_n * u_n$$

RESTRICTED SCHEME

Define

$$e_n := y_n - q_n$$

Perform change of state variable

$$e_n = (d_n + c_n) * u_n$$

$$u_n = g_n * e_n$$

where g_n such that $G(z) = \frac{1}{D(z) + C(z)}$

$$q_n = Q(y_n)$$

$$q_n = Q(y_n)$$

$$y_n - q_n = (d_n + c_n) * u_n$$

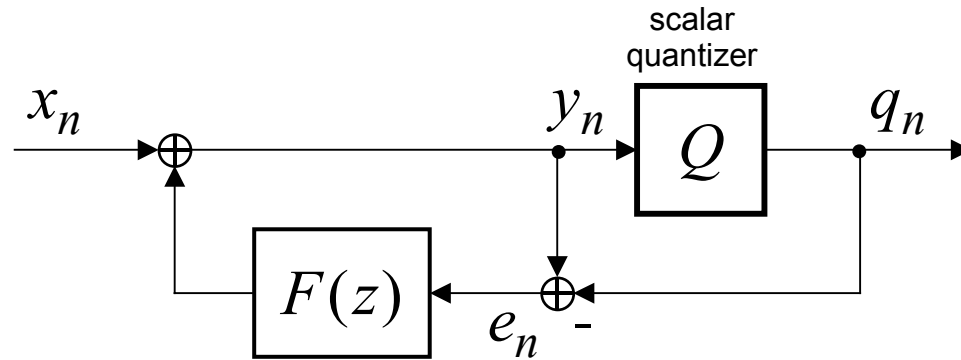
\Leftrightarrow

$$e_n = y_n - q_n$$

$$y_n = x_n + c_n * u_n$$

$$y_n = x_n + \underbrace{c_n * g_n}_{f_n} * e_n$$

ERROR DIFFUSION



$$u_n = g_n * e_n$$

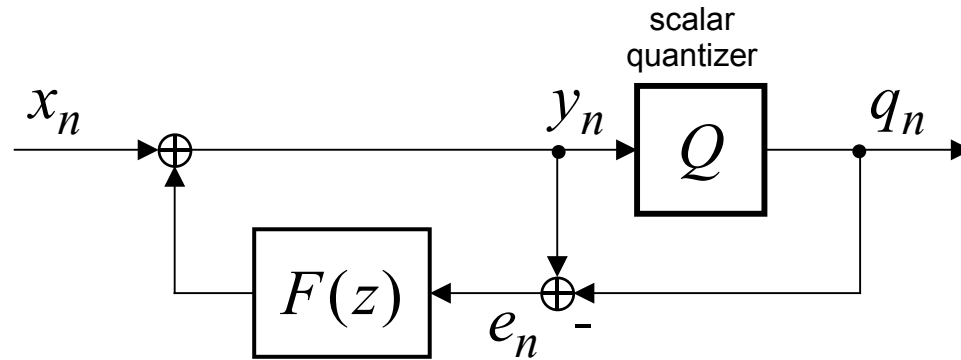
where $G(z) = \frac{1}{D(z) + C(z)}$ and $F(z) = C(z)G(z)$

$$q_n = Q(y_n)$$

$$e_n = y_n - q_n$$

$$y_n = x_n + \underbrace{c_n * g_n}_{f_n} * e_n$$

STABLE ERROR DIFFUSION



$$u_n = g_n * e_n$$

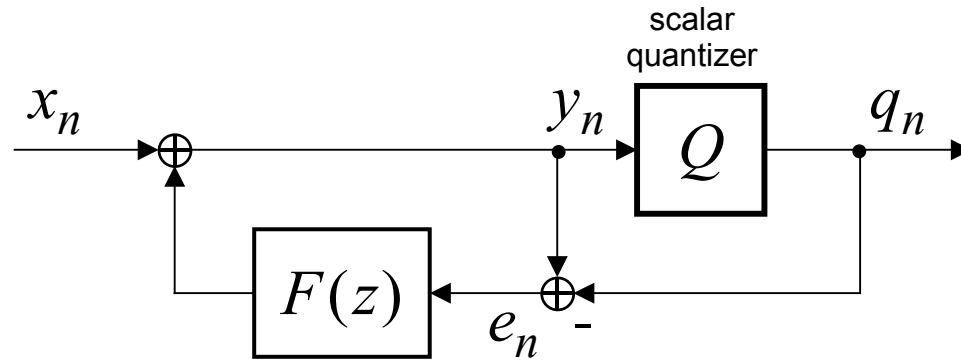
where $G(z) = \frac{1}{D(z) + C(z)}$ and $F(z) = C(z)G(z)$

$$q_n = Q(y_n)$$

$$e_n = y_n - q_n$$

$$\|y\|_\infty \leq \|x\|_\infty + \|f\|_1 \cdot \|e\|_\infty \quad \Leftrightarrow \quad y_n = x_n + f_n * e_n$$

STABLE ERROR DIFFUSION



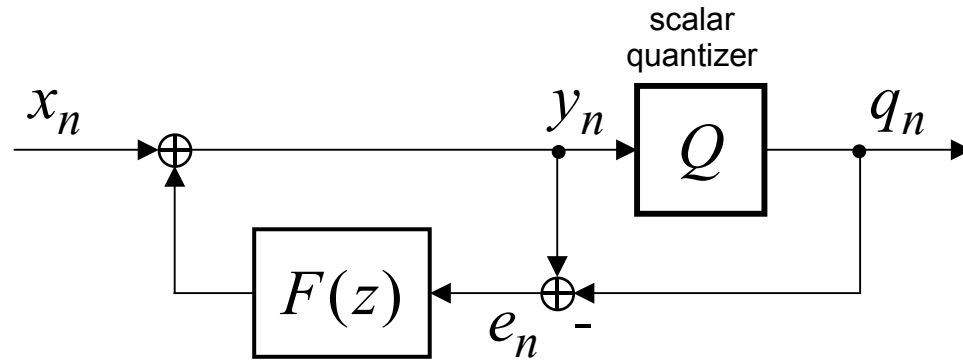
$$u_n = g_n * e_n$$

where $G(z) = \frac{1}{D(z) + C(z)}$ and $F(z) = C(z)G(z)$

- choose $C(z)$ so that both $G(z)$ and $F(z)$ are stable
- choose uniform quantizer Q of step size $\Delta \Rightarrow \|e\|_\infty \leq \frac{\Delta}{2}$

$$\|y\|_\infty \leq \|x\|_\infty + \|f\|_1 \cdot \frac{\Delta}{2} \Rightarrow \text{finite quantizer } Q$$

STABLE 1-BIT ERROR DIFFUSION



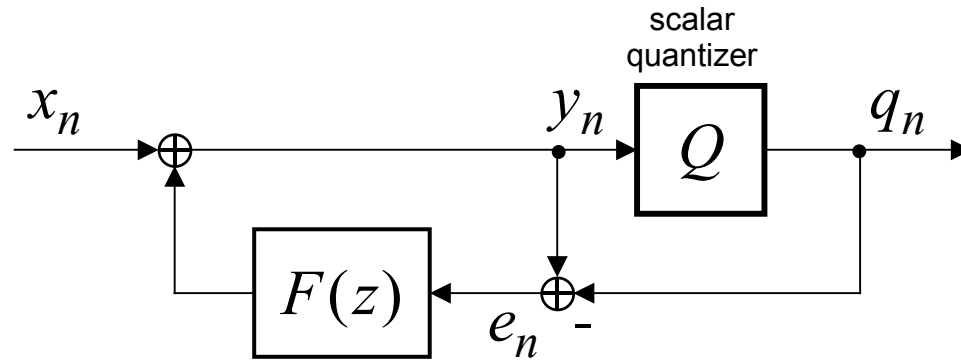
$$u_n = g_n * e_n$$

where $G(z) = \frac{1}{D(z) + C(z)}$ and $F(z) = C(z)G(z)$

- choose $C(z)$ so that both $G(z)$ and $F(z)$ are stable
- choose uniform quantizer Q of step size $\Delta \Rightarrow \|e\|_\infty \leq \frac{\Delta}{2}$
- design f_n so that

$$\|y\|_\infty \leq \|x\|_\infty + \|f\|_1 \cdot \frac{\Delta}{2} \leq \Delta \Rightarrow \text{1-bit quantizer } Q$$

STABLE 1-BIT ERROR DIFFUSION



$$u_n = g_n * e_n$$

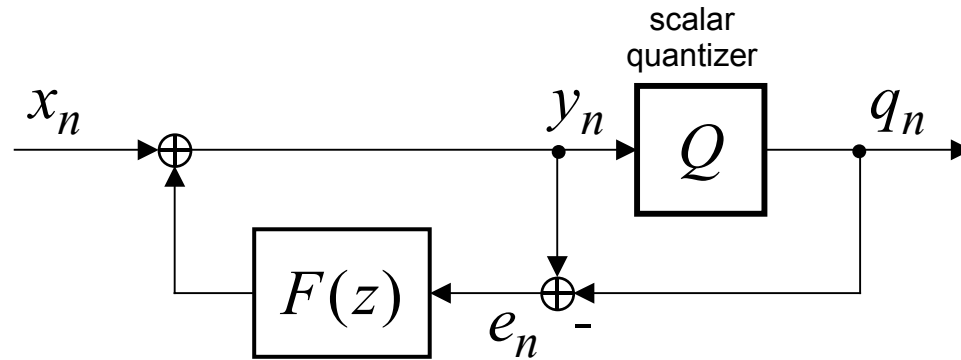
where $G(z) = \frac{1}{D(z) + C(z)}$ and $F(z) = C(z)G(z)$

$\Rightarrow g_n$ causal with $g_0=1$ and $F(z) = 1 - G(z)D(z)$

- design g_n so that

$$\|x\|_\infty + \|f\|_1 \cdot \frac{\Delta}{2} \leq \Delta \Rightarrow \text{1-bit quantizer } Q$$

CASE OF DIFFERENTIATOR



$$D(z) = (1 - z^{-1})^m$$

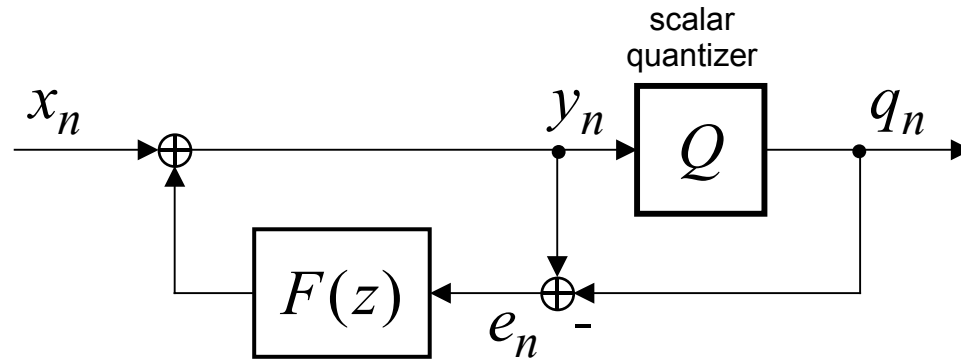
$$F(z) = 1 - G(z)D(z)$$

choose $G(z) = 1$ then $F(z) = 1 - D(z)$

$$\|f_1\| = 2^m - 1$$

$$\|x\|_{\infty} + \|f\|_1 \cdot \frac{\Delta}{2} \leq \Delta \iff \begin{cases} m = 1 \\ \|x\|_{\infty} \leq \frac{\Delta}{2} \end{cases}$$

SINGLE-LOOP SINGLE-BIT CASE



$$D(z) = 1 - z^{-1}$$

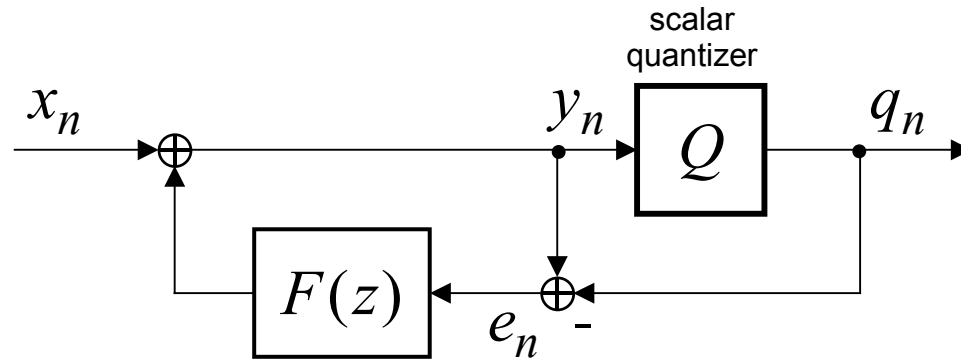
$$F(z) = 1 - G(z)D(z)$$

choose $G(z) = 1$ then $F(z) = z^{-1}$

$$\|f_1\| = 1$$

$$\|x\|_{\infty} + \|f\|_1 \cdot \frac{\Delta}{2} \leq \Delta \iff \begin{cases} m = 1 \\ \|x\|_{\infty} \leq \frac{\Delta}{2} \end{cases}$$

CASE OF DIFFERENTIATOR WITH $m > 1$

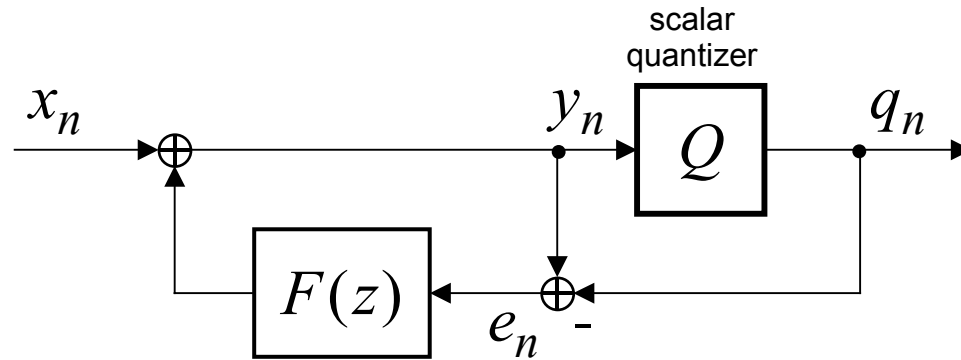


$$D(z) = (1 - z^{-1})^m$$

$$F(z) = 1 - G(z)D(z)$$

$$\|x\|_{\infty} + \|f\|_1 \cdot \frac{\Delta}{2} \leq \Delta \quad \Rightarrow \quad G(z) \neq 1$$

DIFFERENTIATOR WITH FIR $G(z)$



$$D(z) = (1 - z^{-1})^m$$

$$F(z) = 1 - G(z)D(z)$$

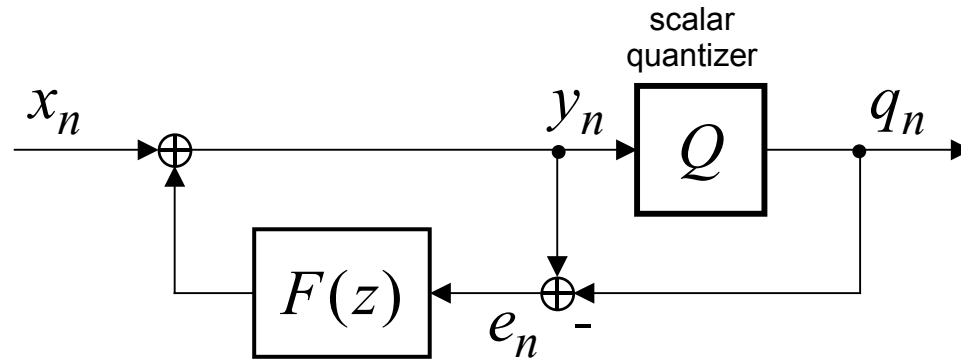
Assumption: $\|x\|_{\infty} \leq \frac{\Delta}{2} - \varepsilon$

There exists FIR filter $G(z)$ such that $\|x\|_{\infty} + \|f\|_1 \cdot \frac{\Delta}{2} \leq \Delta$

$$\text{length}(g_n) \approx 6m^2$$

[S.Gunturk, 2003]

DIFFERENTIATOR WITH IIR $G(z)$



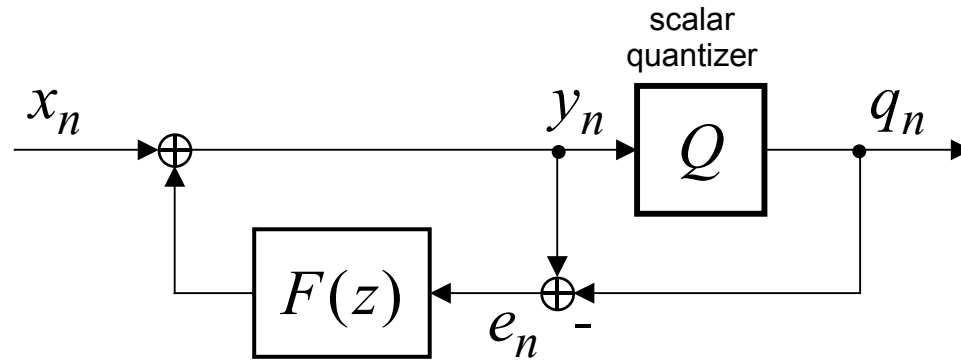
$$D(z) = (1 - z^{-1})^m$$

$$F(z) = 1 - G(z)D(z)$$

with $G(z)$ of the type
$$G(z) = \frac{1}{1 + a_1 z^{-1} + \dots + a_m z^{-m}}$$

- no existing analytical method to guarantee $\|x\|_\infty + \|f\|_1 \cdot \frac{\Delta}{2} \leq \Delta$
- in practice, however, most popular 1-bit scheme !
- empirical design method established by R. Schreier (1993) that leads to (overloaded) stability and high performances

DIFFERENTIATOR WITH IIR $G(z)$: 2nd ORDER



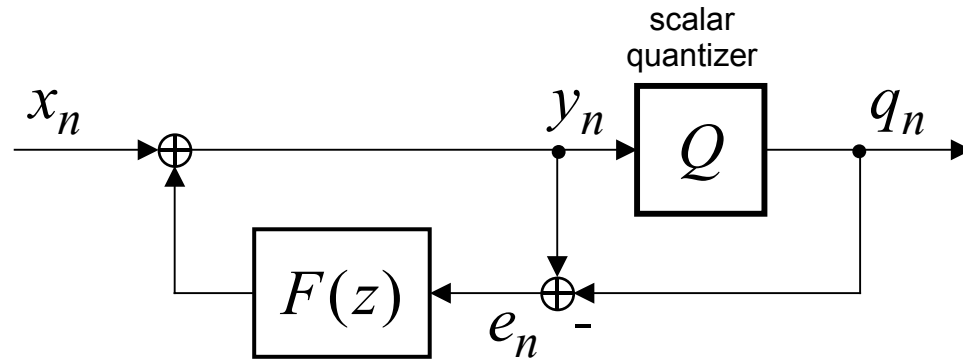
$$D(z) = (1 - z^{-1})^2$$

$$F(z) = 1 - G(z)D(z)$$

with $G(z)$ of the type $G(z) = \frac{1}{1 + a_1 z^{-1} + a_2 z^{-2}}$

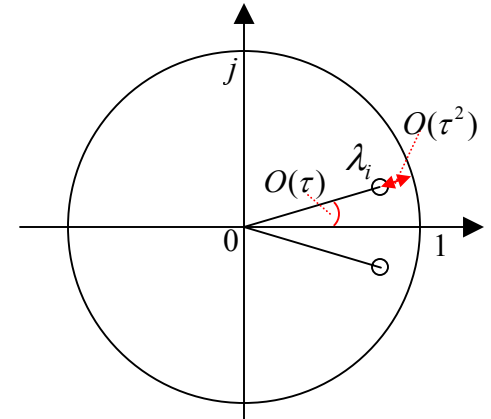
Rigorous proof of (overloaded) stability by [O.Yilmaz, 2002]

LOWPASS $\Sigma\Delta$ WITH IIR $G(z)$



$$D(z) = \prod_{i=1}^m (1 - \lambda_i z^{-1})$$

$$F(z) = 1 - G(z)D(z)$$



with $G(z)$ of the type $G(z) = \frac{1}{1 + a_1 z^{-1} + \dots + a_m z^{-m}}$

- no existing analytical method to guarantee $\|x\|_\infty + \|f\|_1 \cdot \frac{\Delta}{2} \leq \Delta$
- in practice, however, most popular 1-bit scheme !
- empirical design method established by R. Schreier (1993) that leads to (overloaded) stability and high performances

PRINCIPLES OF $\Sigma\Delta$ MODULATION

Choose redundancy operator $D = \{d_{n,k}\}_{n,k \in \mathbb{Z}}$ invertible and

such that $\mathbf{r}_k := \sum_{n \neq k} d_{n,k} \boldsymbol{\varphi}_n$ are “small”

Find quantized sequence $q_n \in \{l_1, l_2, \dots, l_N\}$ so that equation

$$x_n - q_n = \sum_{k \in \mathbb{Z}} d_{n,k} \cdot u_k$$

yields bounded and “small” solution in u_k

Hopefully, $\mathbf{x} - \hat{\mathbf{x}} = \sum_{n \in \mathbb{Z}} u_n \cdot \mathbf{r}_n$ will be “small”

VARIOUS APPLICATIONS

- Bandpass $\Sigma\Delta$ modulation
- Multi-channel $\Sigma\Delta$ modulation
- 2D $\Sigma\Delta$ modulation (image halftoning, time-frequency)
- Finite dimensional space $\Sigma\Delta$ modulation

BANDPASS $\Sigma\Delta$ MODULATION

$$D = \{d_{n-k}\}_{n,k \in \mathbb{Z}} \quad \text{with} \quad d_n = (\delta_n - \delta_{n-1})^{(m)}$$

$$\mathbf{r}_k := \sum_{n \in \mathbb{Z}} d_{n-k} \boldsymbol{\varphi}_n \quad \text{with} \quad \boldsymbol{\varphi}_n = \boldsymbol{\varphi}(t - n\tau)$$



$$\boldsymbol{\varphi}(t) = \cos(\Omega_0 t) \boldsymbol{\psi}(t)$$

↑
lowpass

$$\sum_{k \in \mathbb{Z}} |r_k(t)| \leq \tau^{m-1} \|\boldsymbol{\varphi}^{(m)}\|_{L^1}$$

$$|x(t) - \hat{x}(t)| \leq \|u\|_{\infty} \cdot \sum_{k \in \mathbb{Z}} |r_k(t)|$$

BANDPASS $\Sigma\Delta$ MODULATION

$$D = \{d_{n-k}\}_{n,k \in \mathbb{Z}} \quad \text{with} \quad d_n = \left(\delta_n - e^{j\omega_0} \delta_{n-1}\right)^{(m)} * \left(\delta_n - e^{-j\omega_0} \delta_{n-1}\right)^{(m)}$$

and $\omega_0 = \Omega_0 \tau$

$$\mathbf{r}_k := \sum_{n \in \mathbb{Z}} d_{n-k} \boldsymbol{\varphi}_n \quad \text{with} \quad \boldsymbol{\varphi}_n = \boldsymbol{\varphi}(t - n\tau)$$



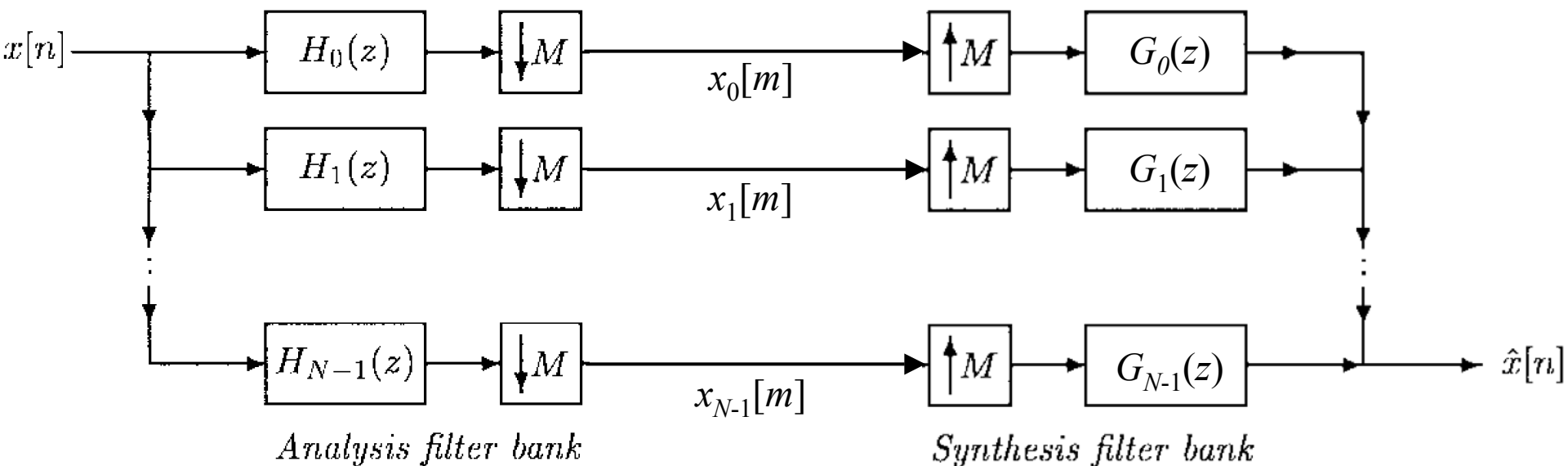
$$\boldsymbol{\varphi}(t) = \cos(\Omega_0 t) \boldsymbol{\psi}(t)$$

↑
lowpass

$$\sum_{k \in \mathbb{Z}} |r_k(t)| \leq (2\tau)^{m-1} \|\boldsymbol{\psi}^{(m)}\|_{L^1}$$

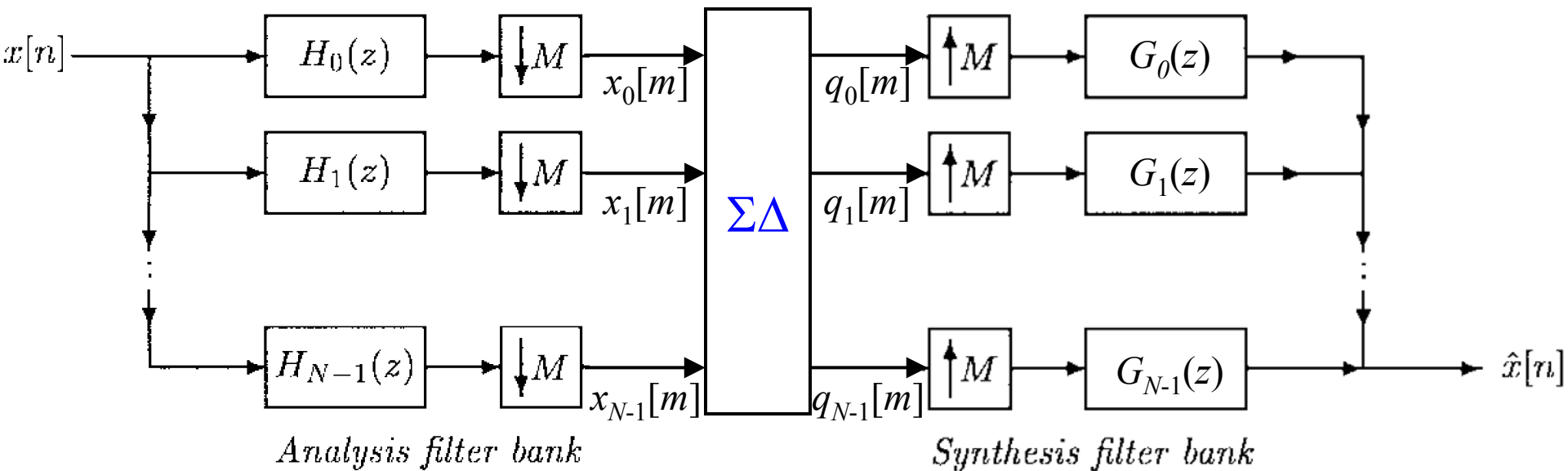
$$|x(t) - \hat{x}(t)| \leq \|u\|_{\infty} \cdot \sum_{k \in \mathbb{Z}} |r_k(t)|$$

MULTI-CHANNEL $\Sigma\Delta$ MODULATION



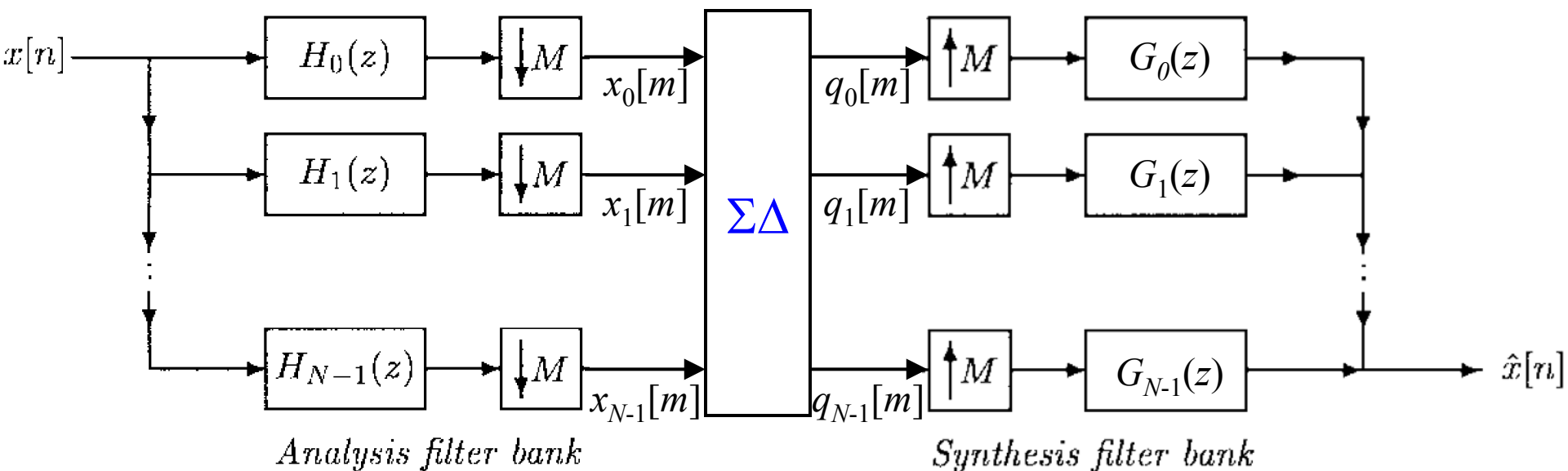
oversampling when $N > M$

MULTI-CHANNEL $\Sigma\Delta$ MODULATION

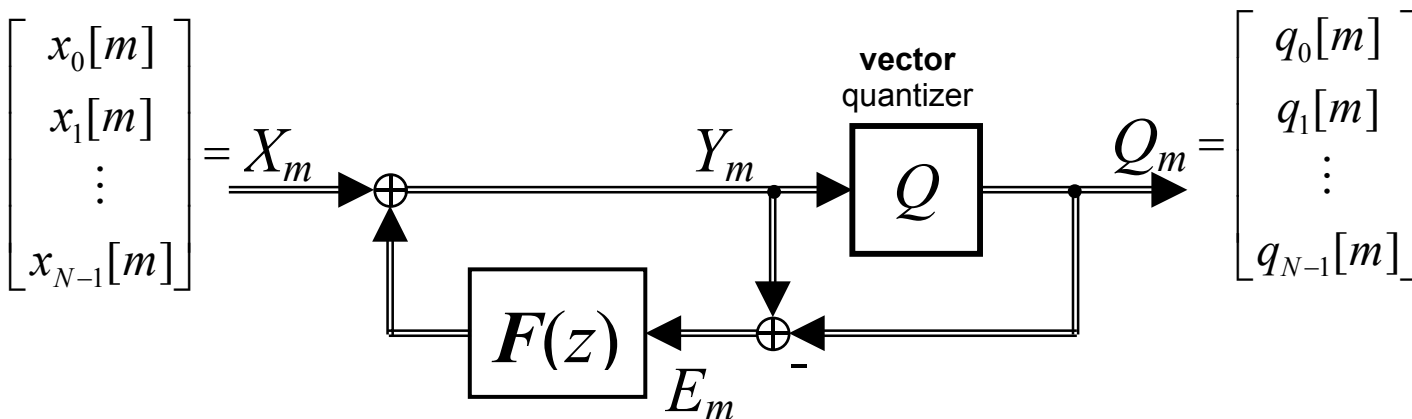


$$N > M$$

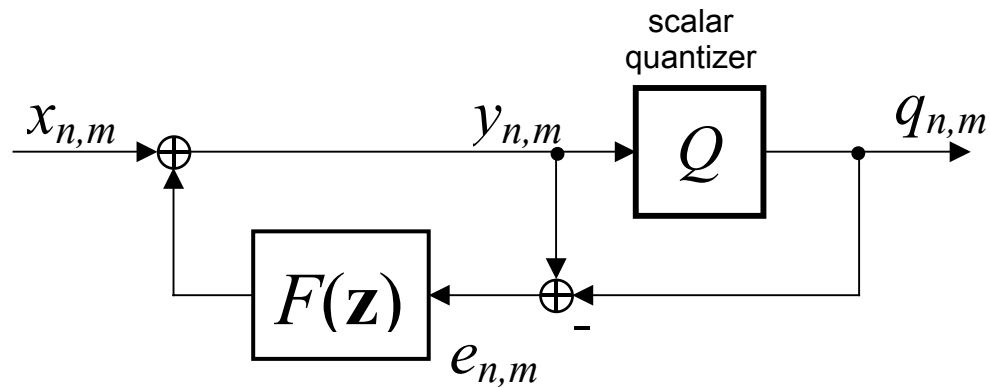
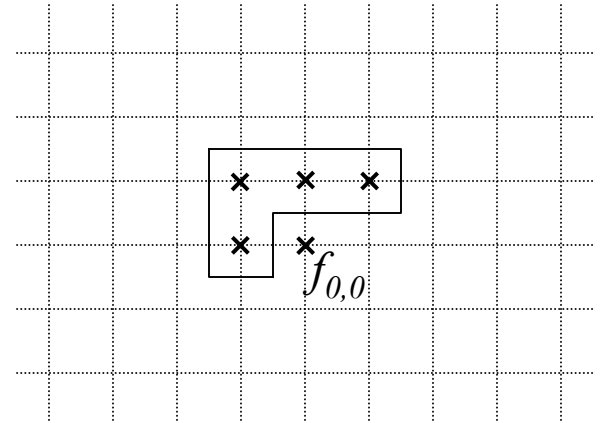
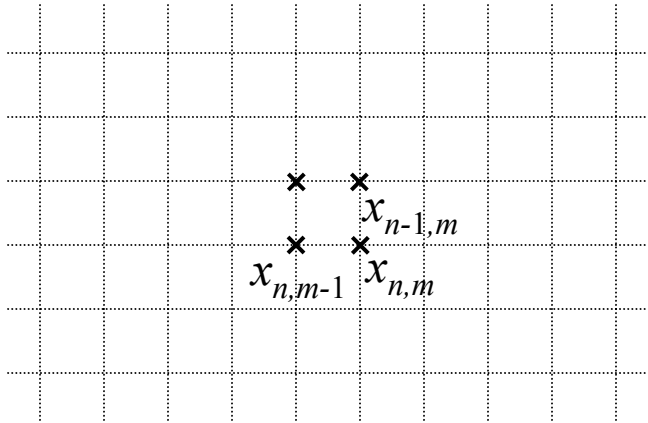
MULTI-CHANNEL $\Sigma\Delta$ MODULATION



$$N > M$$



MULTI-DIMENSIONAL $\Sigma\Delta$ MODULATION



[Image halftoning: C.W.Wu]

[Time-frequency analysis: O.Yilmax, 2004]

INFINITE DIMENSIONAL SPACE

$$\mathbf{x} = \sum_{n \in \mathbb{Z}} x_n \cdot \boldsymbol{\varphi}_n$$

$$\hat{\mathbf{x}} = \sum_{n \in \mathbb{Z}} q_n \cdot \boldsymbol{\varphi}_n$$

$$\mathbf{r}_k := \sum_{n \in \mathbb{Z}} d_{n-k} \cdot \boldsymbol{\varphi}_n$$

$$x_n - q_n = \sum_{k \in \mathbb{Z}} d_{n-k} \cdot u_k$$

$$\mathbf{x} - \hat{\mathbf{x}} = \sum_{n \in \mathbb{Z}} (x_n - q_n) \cdot \boldsymbol{\varphi}_n \quad \Leftrightarrow \quad \mathbf{x} - \hat{\mathbf{x}} = \sum_{k \in \mathbb{Z}} u_k \cdot \mathbf{r}_k$$

FINITE DIMENSIONAL SPACE

$$\mathbf{x} = \sum_{n=1}^N x_n \cdot \boldsymbol{\Phi}_n \quad \hat{\mathbf{x}} = \sum_{n=1}^N q_n \cdot \boldsymbol{\Phi}_n$$

$$\mathbf{r}_k := \sum_{n=1}^N d_{n-k} \cdot \boldsymbol{\Phi}_n, \quad k = 1, \dots, N$$

$$x_n - q_n = \sum_{k=1}^N d_{n-k} \cdot u_k, \quad n = 1, \dots, N$$

$$\mathbf{x} - \hat{\mathbf{x}} = \sum_{n=1}^N (x_n - q_n) \cdot \boldsymbol{\Phi}_n \quad \Leftrightarrow \quad \mathbf{x} - \hat{\mathbf{x}} = \sum_{n=1}^N u_k \cdot \mathbf{r}_k + B$$

↑
boundary
term

[J.Benedetto, O.Yilmaz & A.Powell, 2005]

PRINCIPLES OF $\Sigma\Delta$ MODULATION

Choose redundancy operator $D = \{d_{n,k}\}_{n,k \in \mathbb{Z}}$ invertible and

such that $\mathbf{r}_k := \sum_{n \in \mathbb{Z}} d_{n,k} \boldsymbol{\varphi}_n$ are “small”

Find quantized sequence $q_n \in \{l_1, l_2, \dots, l_N\}$ so that equation

$$x_n - q_n = \sum_{k \in \mathbb{Z}} d_{n,k} \cdot u_k$$

yields bounded and “small” solution in u_k

Hopefully, $\mathbf{x} - \hat{\mathbf{x}} = \sum_{n \in \mathbb{Z}} u_n \cdot \mathbf{r}_n$ will be “small”