

# Finite Word Length Error Analysis based on Basic Formula of Rounding Operation

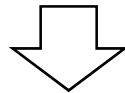
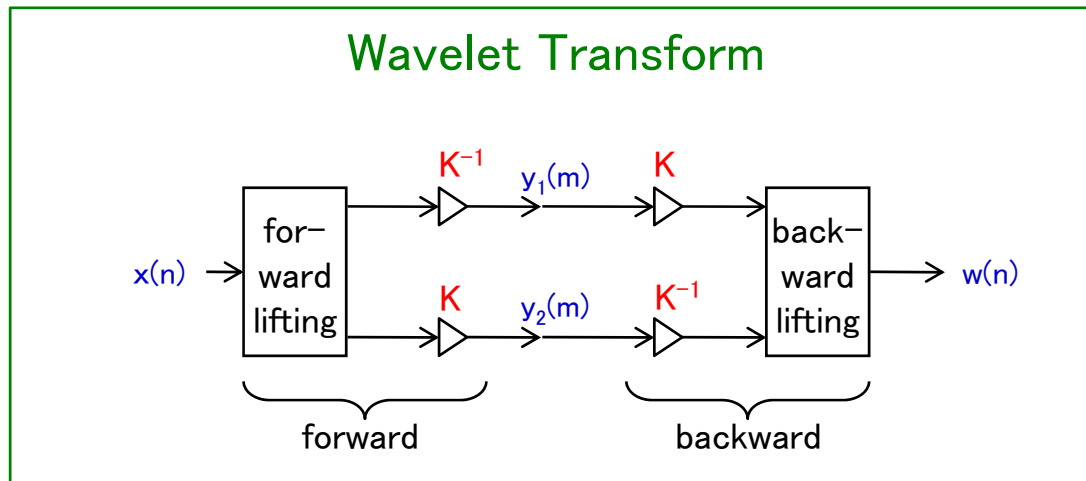
Masahiro IWAHASHI

Nagaoka University of Technology  
Nagaoka-shi, Niigata, 940-2188 Japan

Hitoshi KIYA .

Tokyo Metropolitan University  
Hino-shi, Tokyo, 191-0065 Japan

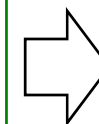
# Wavelet Transform & Perfect Reconstruction (PR)



Perfect Reconstruction

$$w(n) - x(n) = 0$$

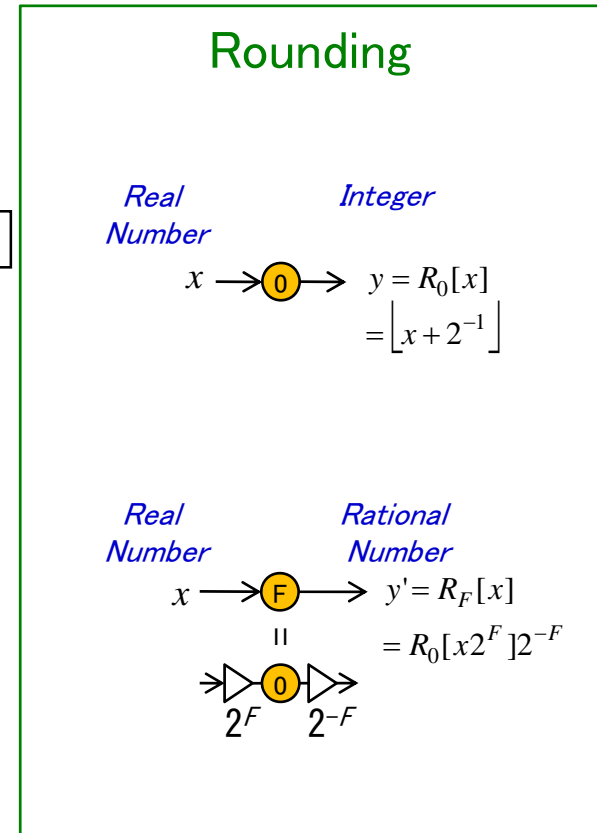
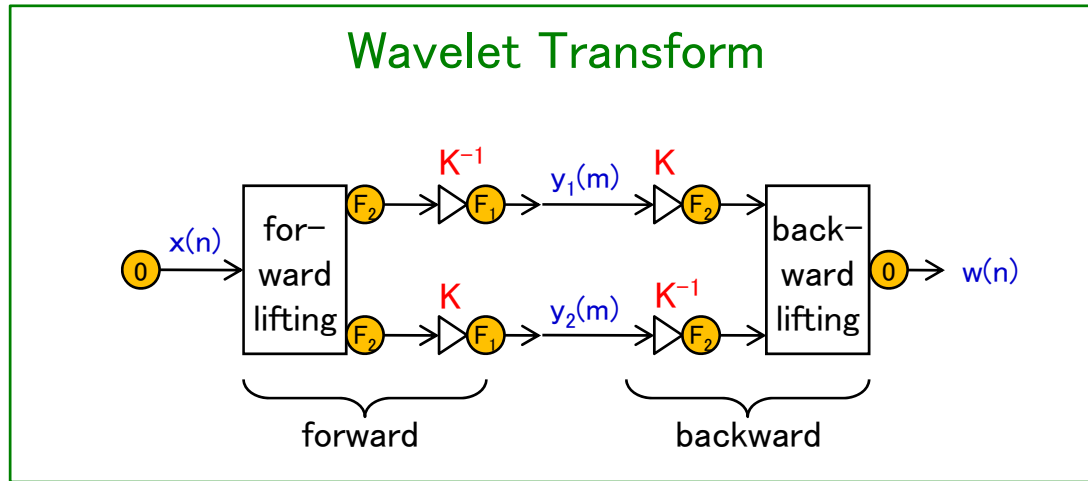
- Signals  $\in$  Real Number
- Coefficients  $\in$  Real Number



assumption  
(design)

# Effect of Rounding

--- Finite Word Length ---



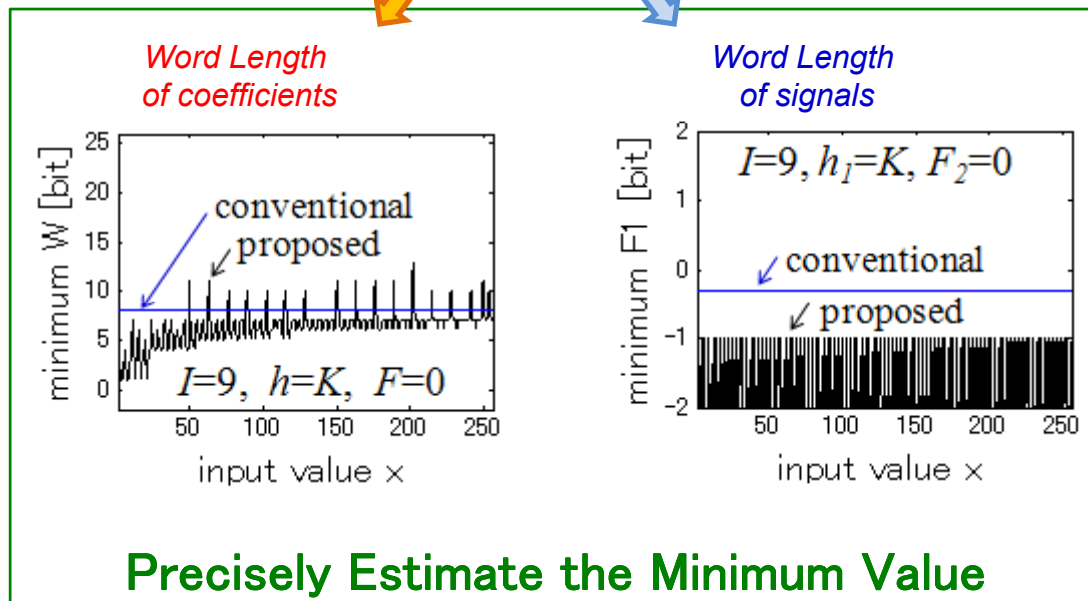
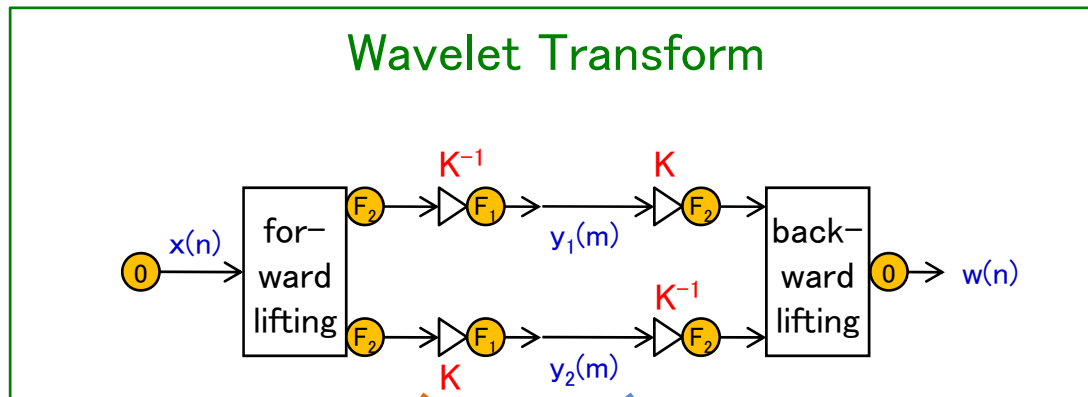
**Perfect Reconstruction**

$w(n) - x(n) \neq 0$

- Signals  $\in$  Integer or Rational Number
- Coefficients  $\in$  Rational Number

$\Rightarrow$  implementation (circuit)

# Find the "Minimum Word Length" which satisfies the Perfect Reconstruction (PR)



**Analysis**

*Conventional ...*

$L^\infty$  norm  
 $\Rightarrow$  precise,  
 $\Rightarrow$  too much, too strict

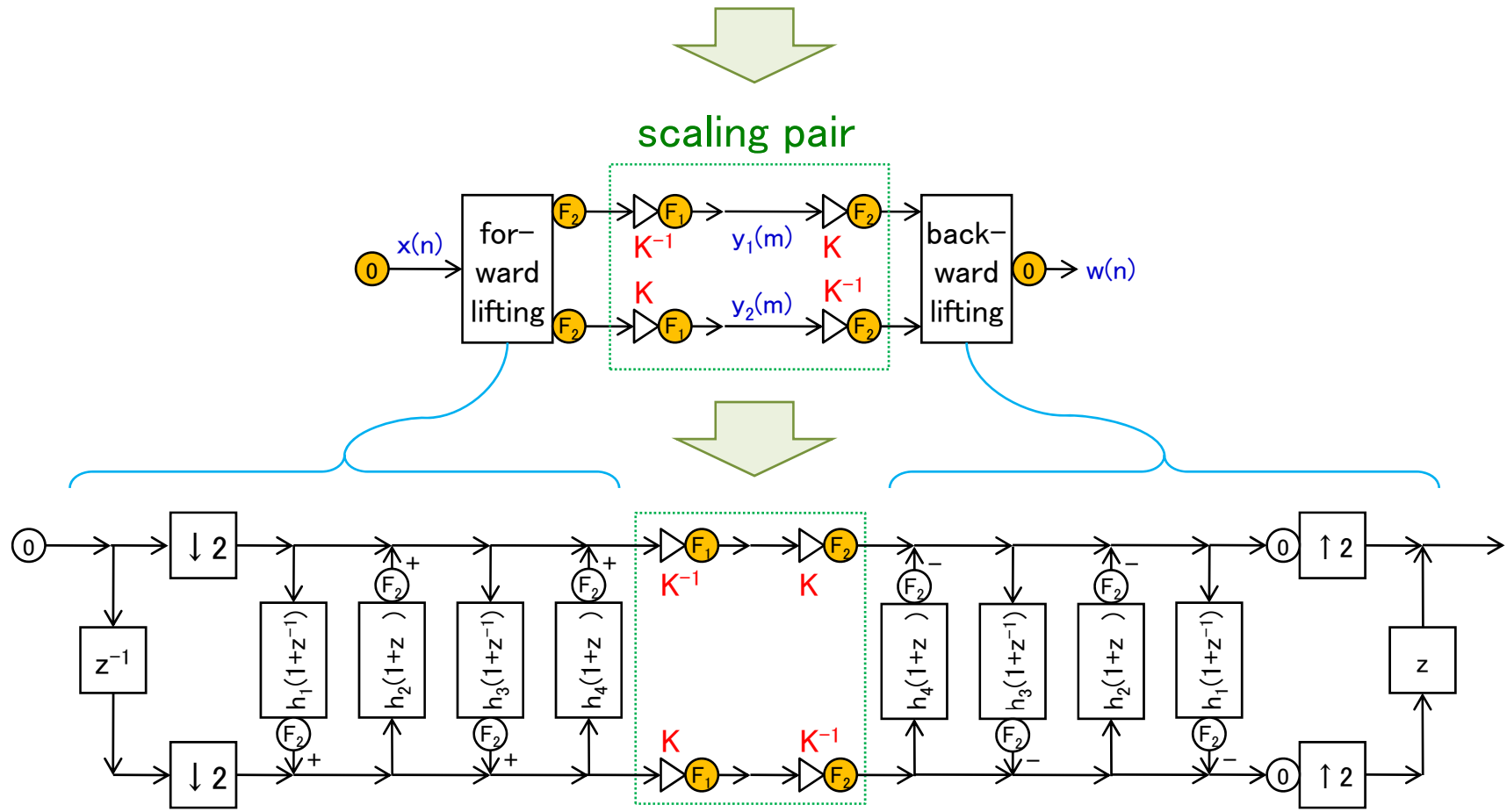
$L^2$  norm  
 $\Rightarrow$  Not Precise  
 $\Rightarrow$  not satisfactory

*We propose ...*

*New analysis*  
 $\Rightarrow$  Precise  
 $\Rightarrow$  exactly same as experimental results

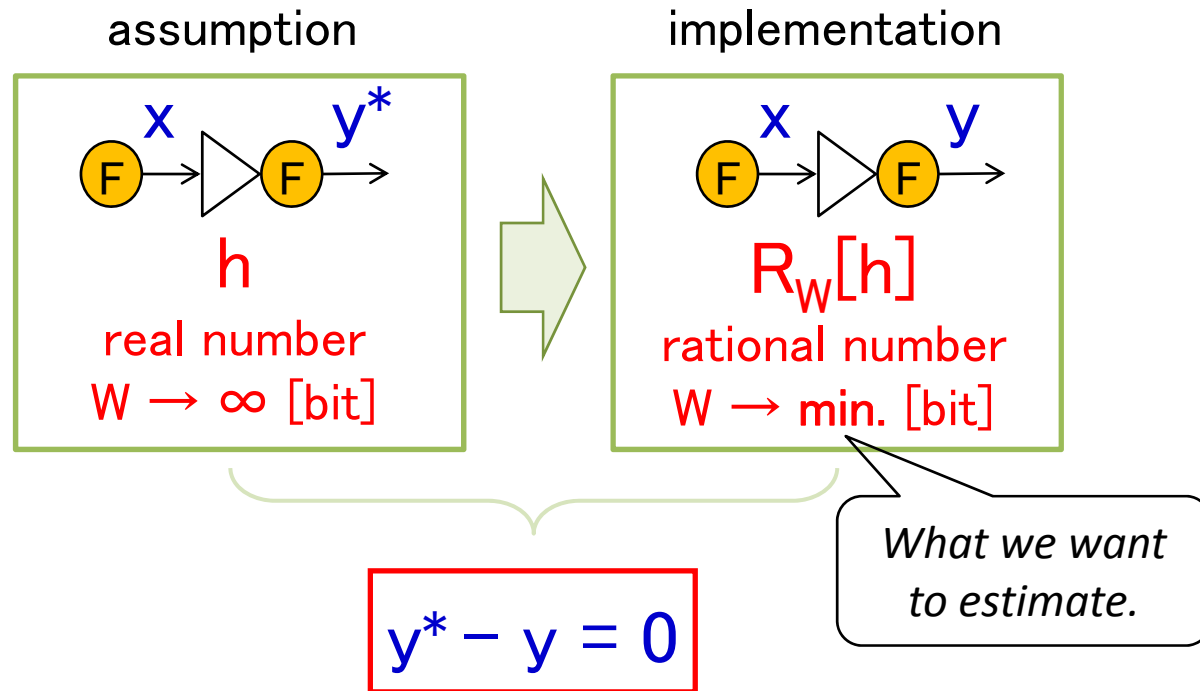
# Conditions on Signals & Coefficients

# Focus on Scaling



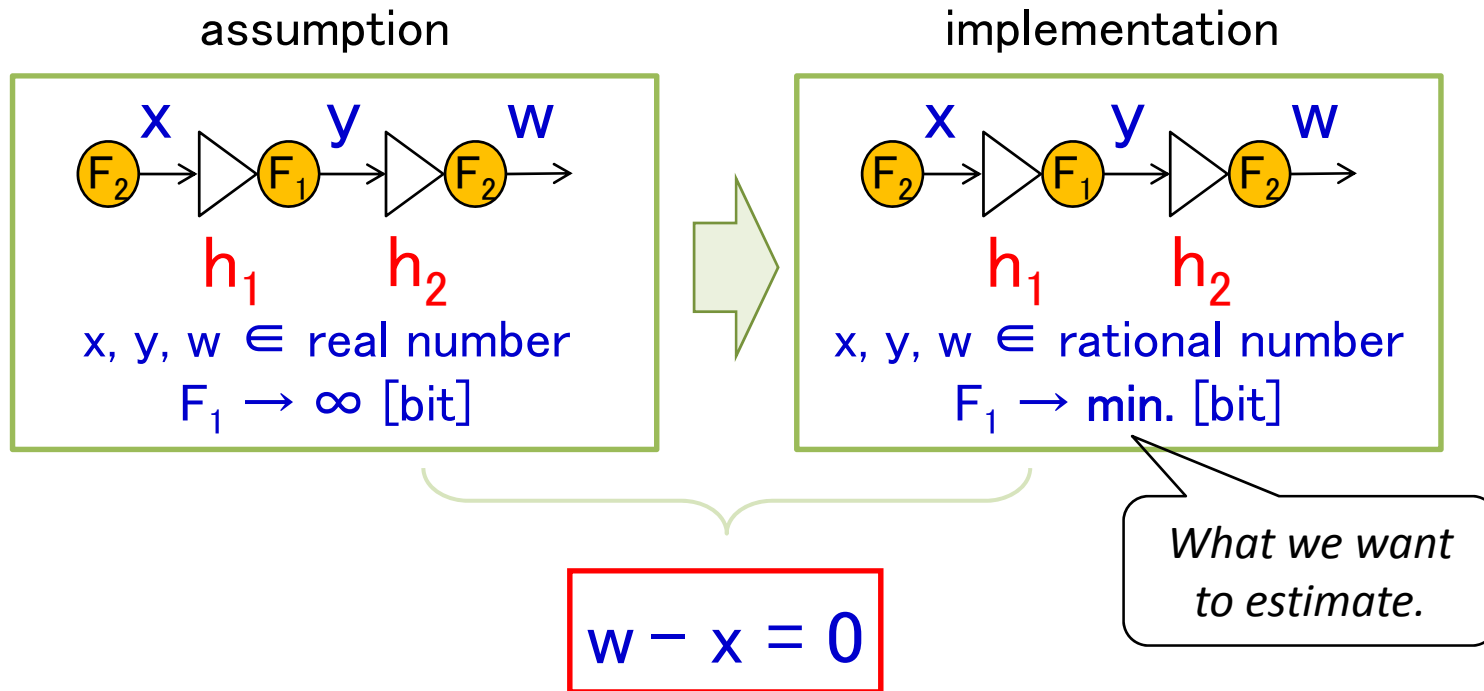
cause of  
 $w(n) - x(n) \neq 0$

# [A] Condition on Coefficient



Mapping Invariant Condition

# [B] Condition on Signals



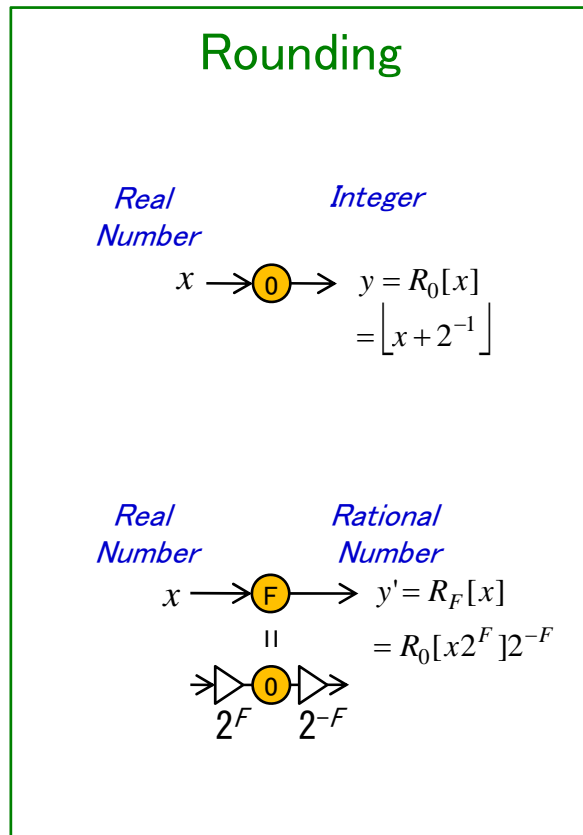
## Lossless Scaling Pair Condition

" $F_2$ " is given and " $h_1 h_2 = 1$ "



# Estimation

# Expression of Rounding



## ● Rounding to Integer

$$\begin{cases} R_0[x] = \lfloor x + 2^{-1} \rfloor = x' - (x' \bmod 1) \\ x' = x + 2^{-1} \end{cases}$$

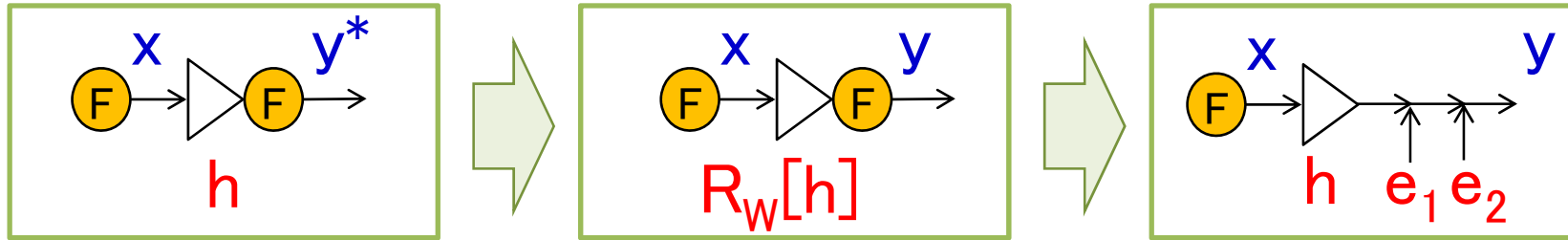
*rounding error*

$$\begin{aligned} \Delta_0[x] &= x - R_0[x] \\ &= \{(x + 2^{-1}) \bmod 1\} - 2^{-1} \end{aligned}$$

## ● Rounding to Rational Number

$$\begin{cases} R_F[x] = R_0[x2^F]2^{-F} \\ \Delta_F[x] = x - R_F[x] = \Delta_0[x2^F]2^{-F} \end{cases}$$

# [A] Condition on Coefficient



## Conventional Analysis

$L^\infty$  norm

⇒ precise,

⇒ too much, too strict

$L^2$  norm

⇒ Not Precise

⇒ not appropriate

$$\max|e_1| = 2^{-1-W} 2^{-1+I} < \max|e_2| = 2^{-1-F}$$

$$\sigma_{e_1}^2 = \frac{1}{12} (2^{-1-W} 2^{-1+I} \cdot 2)^2 \ll \sigma_{e_2}^2 = \frac{1}{12} (2^{-1-F} \cdot 2)^2$$

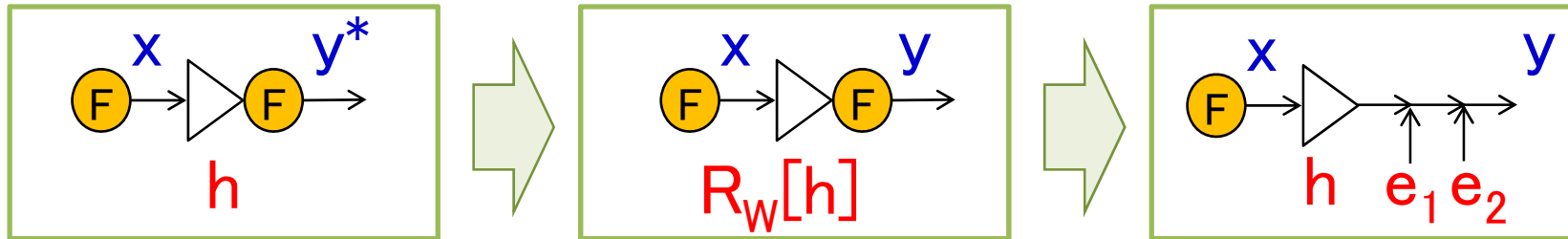


$$W > F + I - 1$$

$$\begin{cases} y = hx + e_1 + e_2 \\ e_1 = -\Delta_W[h]x \\ e_2 = -\Delta_F[W[h]x] \end{cases}$$

estimated.

# [A] Condition on Coefficient



$$\begin{cases} y = hx + e_1 + e_2 \\ e_1 = -\Delta_W[h]x \\ e_2 = -\Delta_F[W[h]x] \end{cases}$$

## Our Analysis

*New analysis*

*⇒ Precise*

*⇒ exactly same as*

*experimental results*

$$E_m = 0 \quad \text{for} \quad \begin{cases} E_m = R_F[R_W[h]x] - R_F[hx] \\ R_W[h] = h - \Delta_W[h] \end{cases}$$



$$|\Delta_W[h]x - \Delta_F[hx]| < 2^{-1-F} \quad \text{estimated.}$$

# [A] proof

$$E_m = 0 \quad \text{for} \quad \begin{cases} E_m = R_F[R_W[h]x] - R_F[hx] \\ R_W[h] = h - \Delta_W[h] \end{cases}$$

*proof*

$$\begin{aligned} & R_F[R_W[h]x] - R_F[hx] \\ &= R_F[hx - \Delta_W[h]x] - R_F[hx] \\ &= R_F[\Delta_F[hx] - \Delta_W[h]x] + R_F[hx] - R_F[hx] \\ &= R_F[\Delta_F[hx] - \Delta_W[h]x] = 0 \end{aligned}$$

$$\therefore |\Delta_F[hx] - \Delta_W[h]x| < 2^{-1-F}$$

**Q.E.D.**

## Addition Formula

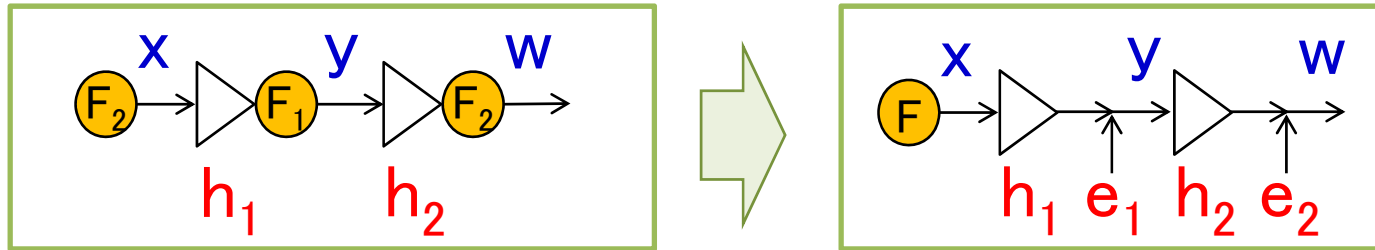
$$R_F[x + y] = R_F[x] + R_F[y + \Delta_F[x]]$$

$x \in \mathbf{R}, y \in \mathbf{R}$

## Properties

$$\begin{cases} R_0[x] = 0 \Leftrightarrow x \in [-2^{-1}, 2^{-1}) \\ R_F[x] = 0 \Leftrightarrow x \in [-2^{-1-F}, 2^{-1-F}) \end{cases}$$

# [B] Condition on Signals



$$\begin{cases} w = h_2 y + e_2 \\ y = h_1 x + e_1 \end{cases}, \quad \begin{cases} e_1 = -\Delta_{F_1} [h_1 x] \\ e_2 = -\Delta_{F_2} [h_2 y] \end{cases}$$

## Conventional Analysis

$L^\infty$  norm

$\Rightarrow$  precise,  
 $\Rightarrow$  too much, too strict

$L^2$  norm

$\Rightarrow$  Not Precise  
 $\Rightarrow$  not appropriate

$$\max |h_2 e_1| = |h_2| 2^{-1+F_1} < \max |e_2| = 2^{-1-F_2}$$

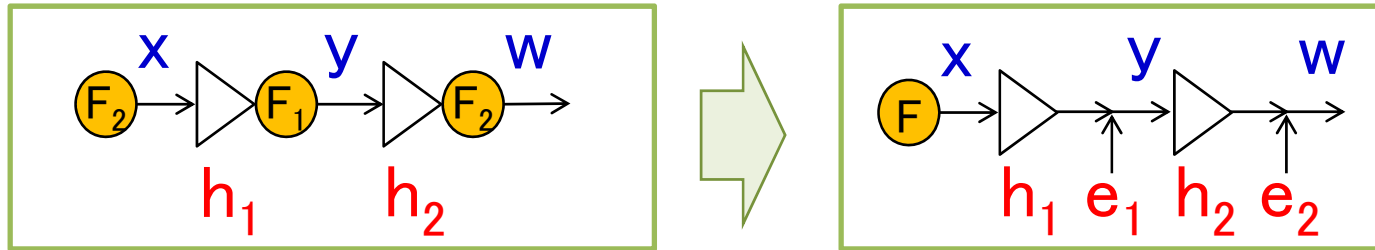
$$\sigma_{h_2 e_1}^2 = \frac{1}{12} (h_2 2^{-1+F_1} \cdot 2)^2 \ll \sigma_{e_2}^2 = \frac{1}{12} (2^{-1-F_2} \cdot 2)^2$$



$$F_1 > F_2 + \log_2 |h_2|$$

estimated.

# [B] Condition on Signals



$$\begin{cases} w = h_2 y + e_2 \\ y = h_1 x + e_1 \end{cases}, \quad \begin{cases} e_1 = -\Delta_{F_1} [h_1 x] \\ e_2 = -\Delta_{F_2} [h_2 y] \end{cases}$$

## Our Analysis

*New analysis*

*⇒ Precise*

*⇒ exactly same as experimental results*

$$E_p = 0 \quad \text{for} \quad \begin{cases} E_p = R_{F_2} [h_2 R_{F_1} [h_1 x]] - x \\ h_1 h_2 = 1 \end{cases}$$



$$\left| h_2 \Delta_{F_1} [h_1 x] \right| < 2^{-1-F_2} \quad \text{estimated.}$$

# [B] proof

$$\begin{aligned} & R_{F_2} [xR_{F_1} [y]] \\ &= R_{F_2} [xy] + R_{F_2} [-x\Delta_{F_1} [y] + \Delta_{F_2} [xy]] \end{aligned}$$



## Multiplication Formula

$$\begin{aligned} & R_F [xy] \\ &= R_F [xR_F [y]] + R_F [x\Delta_F [y] + \Delta_F [xR_F [y]]] \\ & x \in \mathbf{R}, y \in \mathbf{R} \end{aligned}$$

*proof*

$$\begin{aligned} & R_{F_2} [xR_{F_1} [y]] \\ &= R_0 [-x\Delta_{F_1} [y]2^{F_2} + xy2^{F_2}]2^{-F_2} \\ &= R_0 [-x\Delta_{F_1} [y]2^{F_2} + \Delta_0 [xy2^{F_2}]]2^{-F_2} + R_{F_2} [xy] \\ &= R_{F_2} [-x\Delta_{F_1} [y] + \Delta_{F_2} [xy]] + R_{F_2} [xy] \end{aligned}$$



**Q.E.D.**

$$R_{F_2} [h_2\Delta_{F_1} [h_1x]] = 0$$

$$|h_2\Delta_{F_1} [h_1x]| < 2^{-1-F_2}$$



## Properties

$$\begin{cases} R_0 [x] = 0 \Leftrightarrow x \in [-2^{-1}, 2^{-1}) \\ R_F [x] = 0 \Leftrightarrow x \in [-2^{-1-F}, 2^{-1-F}) \end{cases}$$



# *Proof of the Formulas*

## *Addition Formula*

$$R_F[x + y] = R_F[x] + R_F[y + \Delta_F[x]]$$

$x \in \mathbf{R}, y \in \mathbf{R}$

*proof*

$$\begin{aligned} R_F[x + y] &= R_F[R_F[x] + \Delta_F[x] + y] \\ &= R_0[R_0[x2^F] + \Delta_F[x]2^F + y2^F]2^{-F} \\ &= R_F[x] + R_F[\Delta_F[x] + y] \end{aligned}$$

*Q.E.D.*

## *Multiplication Formula*

$$R_F[xy] = R_F[xR_F[y]] + R_F[x\Delta_F[y] + \Delta_F[xR_F[y]]]$$

$x \in \mathbf{R}, y \in \mathbf{R}$

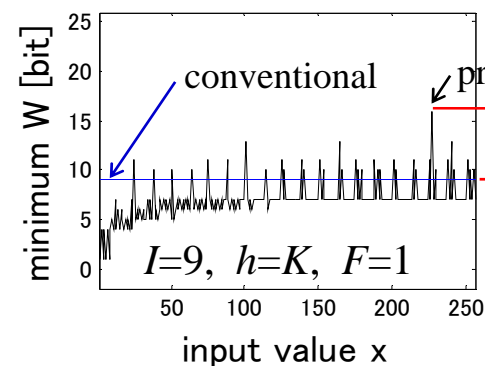
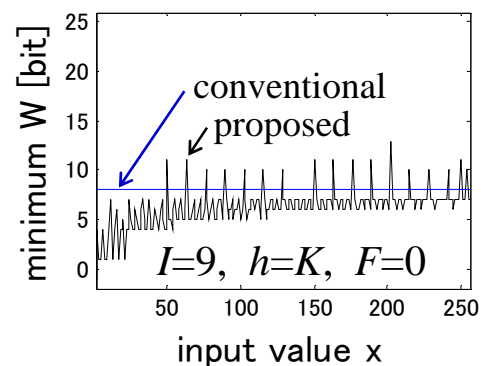
*proof*

$$\begin{aligned} R_F[xy] &= R_F[x\Delta_F[y] + xR_F[y]] \\ &= R_F[x\Delta_F[y] + \Delta_F[xR_F[y]] + R_F[xR_F[y]] \\ &= R_F[x\Delta_F[y] + \Delta_F[xR_F[y]]] + R_F[xR_F[y]] \end{aligned}$$

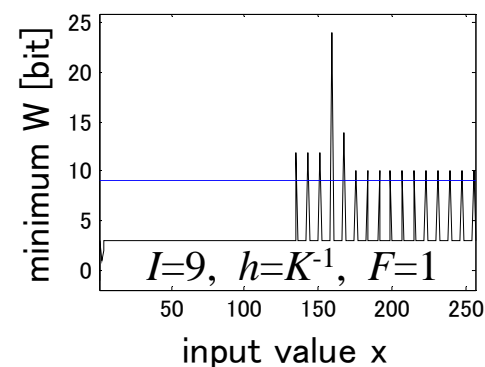
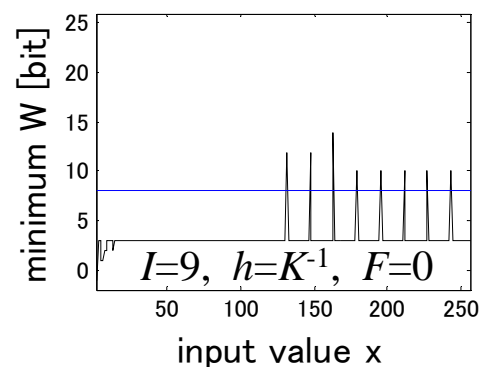
*Q.E.D.*

# Results

# Min. Word Length of Coefficients

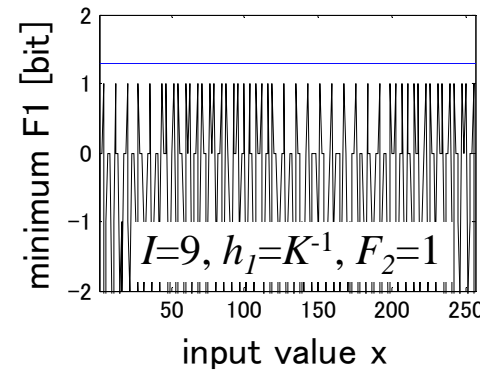
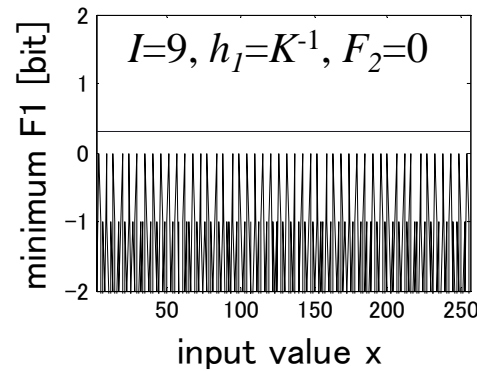
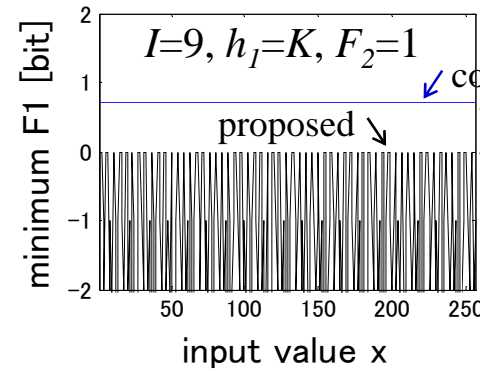
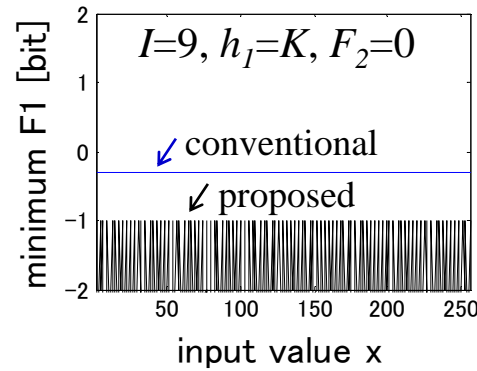


*Result of the  
conventional  
analysis is  
not satisfactory.*



$$W_{\min} = \{W \in \mathbf{Z} \mid E_m = 0, \forall W \geq W_{\min}\}$$

# Min. Word Length of Signals



*Result of the conventional analysis is too strict.*

$$F_{1,\min} = \{F_1 \in \mathbf{Z} \mid E_p = 0, \forall F_1 \geq F_{1,\min}\}$$

# Conclusion

# Conclusion Remarks

In this report, we mathematically derived formulas of the rounding operation

to analyze the lossless condition on word length of signals and coefficients.

It was confirmed that our equations can determine the minimum word length for each input value more precisely than the conventional analysis.