

# **ABCD : Arbitrary Bitwise Coefficient for De-Quantization**

**Woo Kyoung Han, Byeong Hun Lee, Sang Hyun Park, Kyong Hwan Jin**

**TUE-PM-167**

# Visual Demonstration

Sintel **2-bit** → 8-bit



2-bit Input

ABCD (Ours)

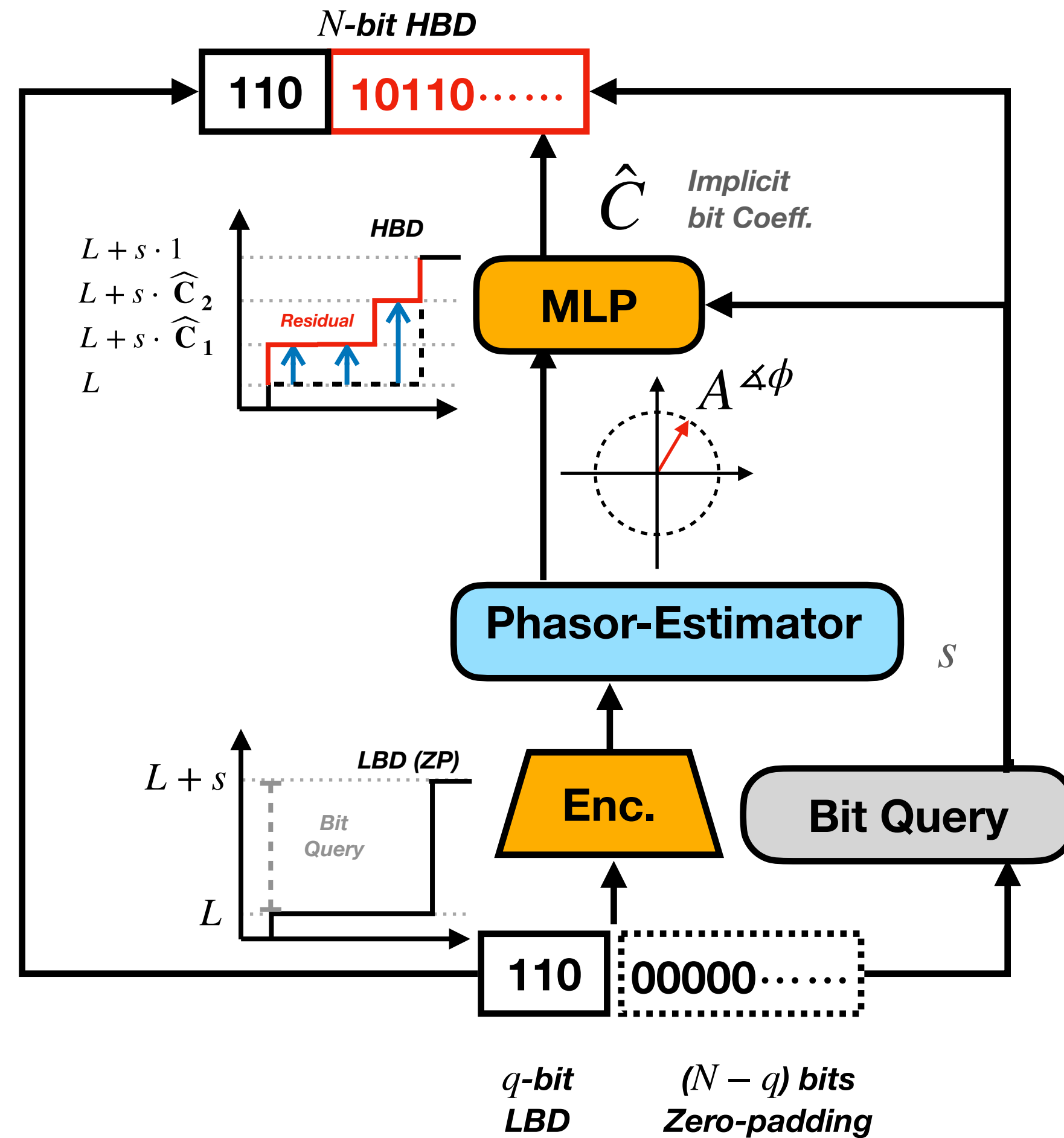
GT

Blurred Details



False Contour

# ABCD



- BDE algorithm using INR  
**With Bitwise Coefficients**
- Phasor estimator  
**To Relieve Spectral bias**
- State-of-the-Art Performance  
**with Single training**

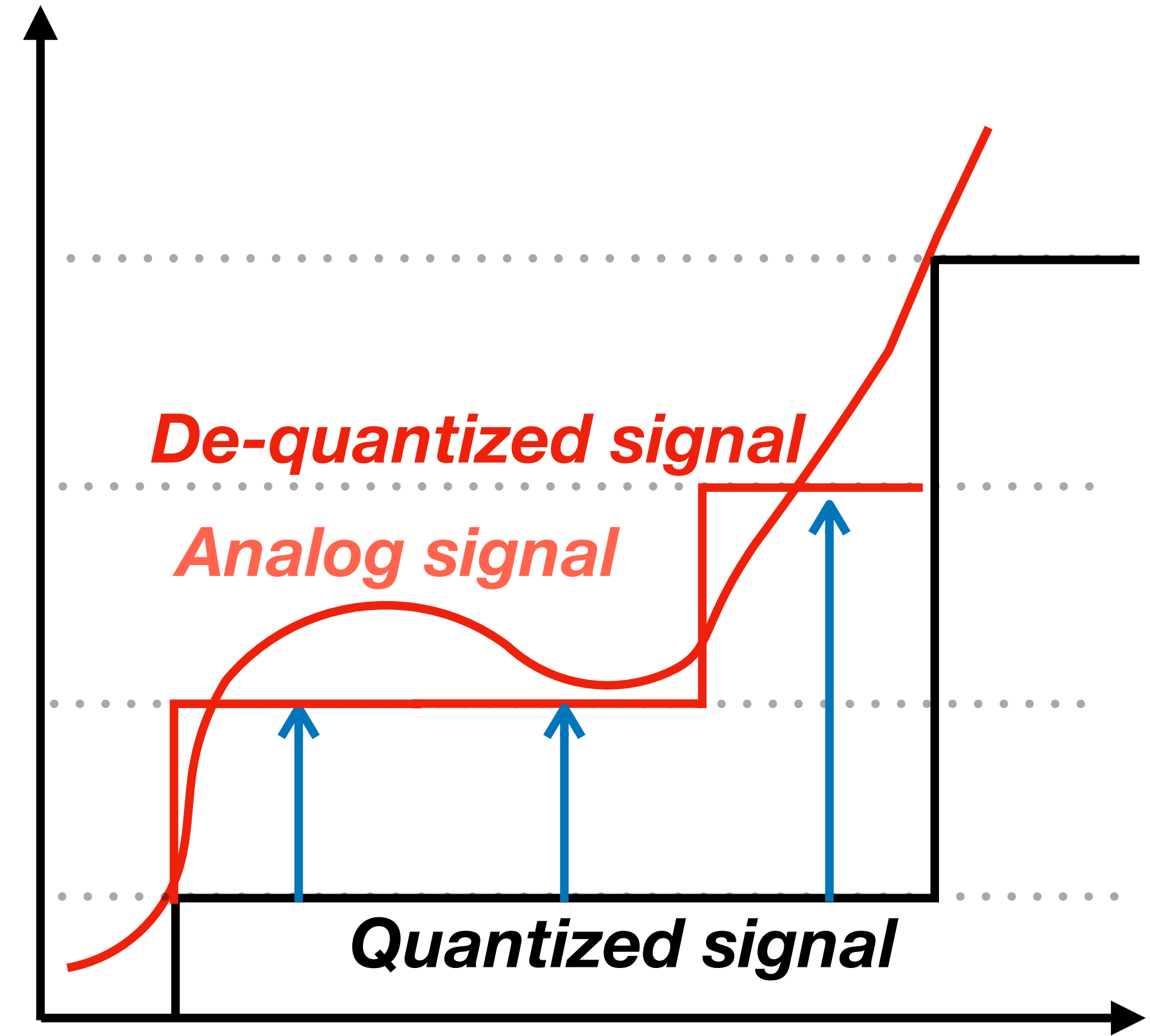
# Bit Depth Expansion

Quantization & De-quantization

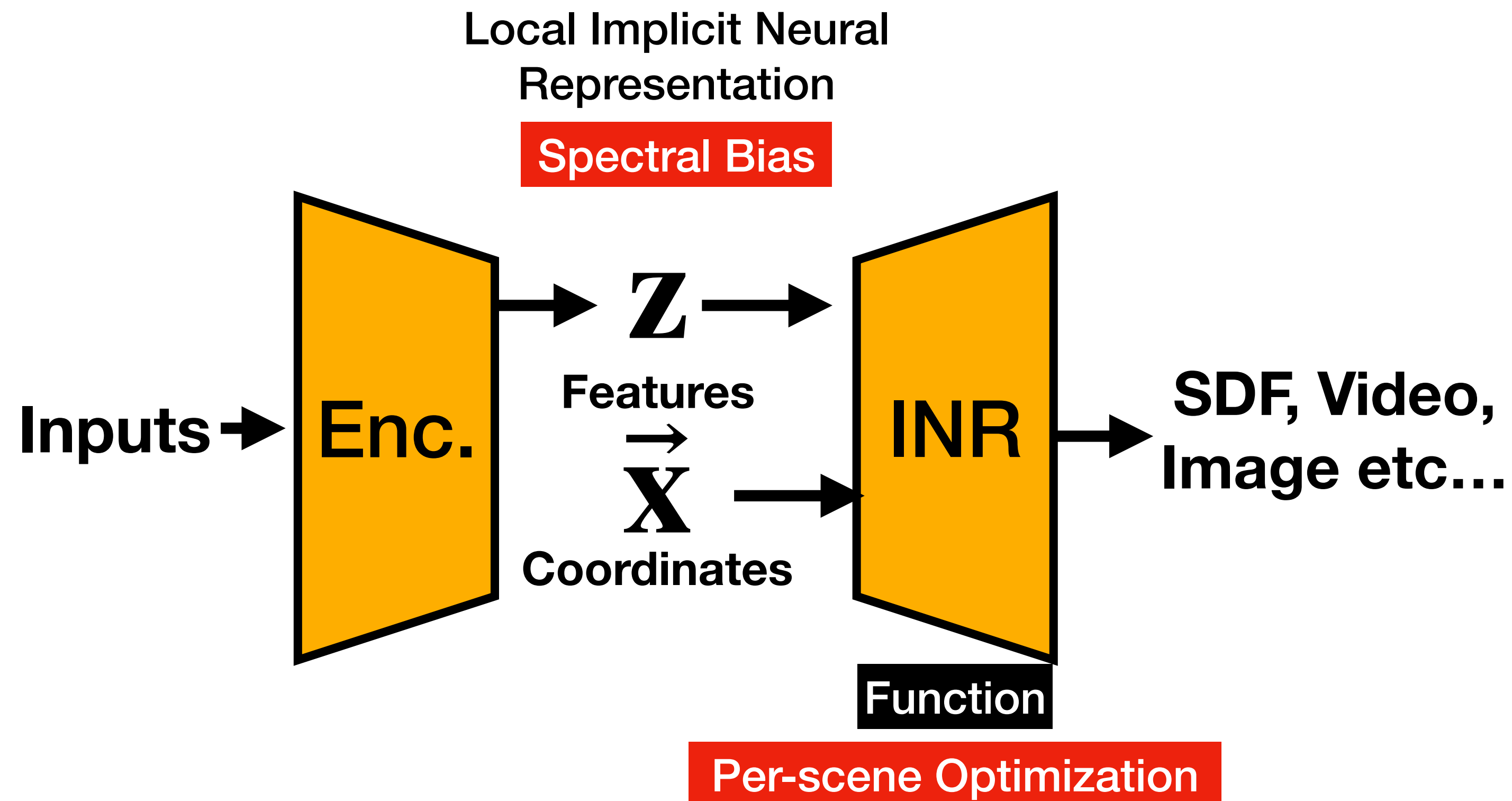


*High Bit Depth (HBD)*

110	100111011
-----	-----------



# Implicit Neural Representation



- How to define a coordinate and corresponding signal to BDE?
- How to mitigate the spectral bias?

# Formulation

## Bit-wise Coefficient

$$\begin{aligned}
 a &= \sum_{i=-\infty}^{\infty} b_i 2^i && (b_i \in \{0,1\}) \\
 \text{positive real number} &&& \\
 &= \sum_{i=L+1}^{\infty} b_i 2^i + \sum_{j=-\infty}^L b_j 2^j && (* \sum_{k=0}^{\infty} a_0 r^k = \frac{a_0}{1-r} \quad (|r| < 1)) \\
 &= \sum_{i=L+1}^{\infty} b_i 2^i + C \cdot 2^{L+1}, && C \in [0,1]
 \end{aligned}$$

# Formulation

## Bit-wise Basis

$$\mathbb{F}_2 := (\{0,1\}, \oplus, \cdot)$$

( $\oplus$  : XOR operation,  $\cdot$  : multiplication)

$$\{0,1\}^N \longrightarrow \exists e_q$$

Orthonormal basis

$$e_q \sim 2^{N-q}$$

(As digital number)

• ex)  $e_4 =$  00010000  $\sim 2^{8-4}$

# Formulation

## Image Quantization

$$\mathbf{I}_q = \left\lfloor \frac{\mathbf{I}_N}{2^{N-q}} \right\rfloor 2^{N-q}$$

$$\mathbf{I}_N = \mathbf{I}_q + \mathbf{R}$$

$$= \underbrace{b_1 b_2 b_3}_{\mathbf{I}_q} \underbrace{b_4 b_5 b_6 b_7 b_8}_{\mathbf{R}}$$

## Bit-wise Coefficients

$$a = \sum_{i=L+1}^{\infty} b_i 2^i + \mathbf{C} \cdot 2^{L+1}$$

## Bit-wise basis

$$e_q \sim 2^{N-q}$$

## Conclusion

$$\mathbf{I}_N = \underbrace{\sum_{i=N-q}^{N-1} 2^i \cdot \mathbf{B}_i}_{\mathbf{I}_q} + \underbrace{2^{N-q} \cdot \mathbf{C}}_{\mathbf{R}} \rightarrow f_{\theta}(e_q, \mathbf{z})$$

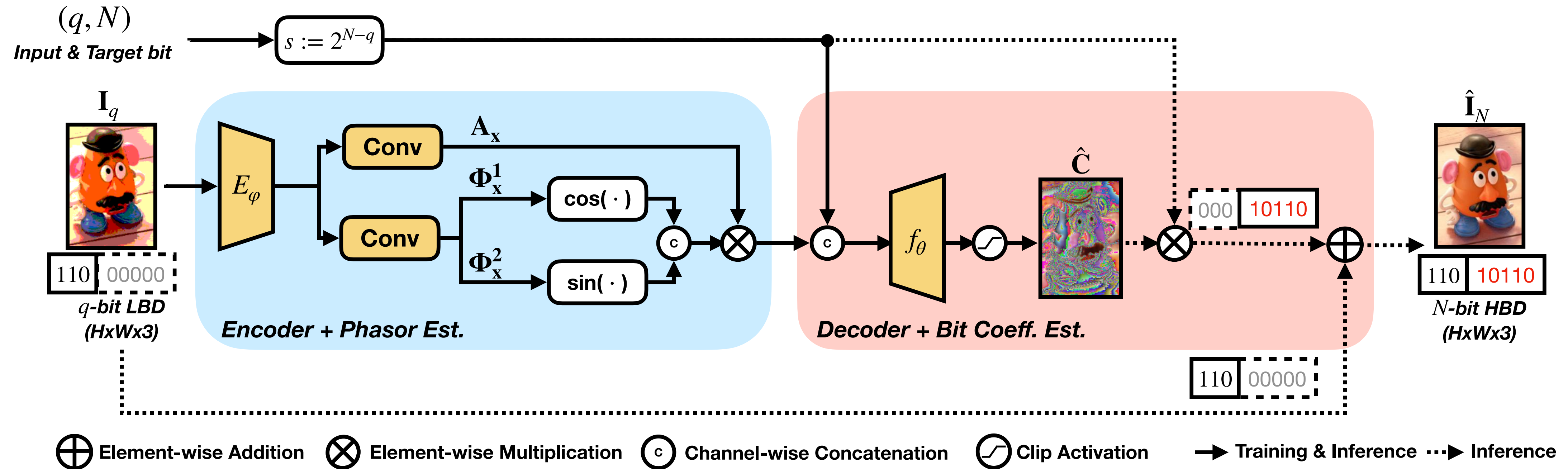


# Formulation

## Phasor estimator

$$\underbrace{\mathbf{A}_j \odot \begin{bmatrix} \cos(\pi(\mathbf{F}_j \cdot \delta + h_p(\hat{\mathbf{c}}))) \\ \sin(\pi(\mathbf{F}_j \cdot \delta + h_p(\hat{\mathbf{c}}))) \end{bmatrix}}_{\text{LTE}} \Big|_{\vec{\delta} = \vec{\mathbf{0}}} \longrightarrow \underbrace{\begin{bmatrix} \mathbf{A}_x^1 \\ \mathbf{A}_x^2 \end{bmatrix} \odot \begin{bmatrix} \cos(\pi\Phi_x^1) \\ \sin(\pi\Phi_x^2) \end{bmatrix}}_{\text{ABCD}}$$

# ABCD Structure



- Encoder ( $E_\varphi$ ) : EDSR (17'), RDN (18'), SwinIR (21')
- Decoder ( $f_\theta$ ) : 4-Layer Multi-layer perceptron (MLP)

# Quantitative comparison

## PSNR(dB) / SSIM

Vaild Method	Sintel						MIT-Adobe FiveK			
	4 >> 8	4 >> 12	4 >> 16	6 >> 12	6 >> 16	8 >> 16	3 >> 16	4 >> 16	5 >> 16	6 >> 16
Input(zp)	29.16	28.78	28.77	40.90	40.81	52.85	22.90	28.86	34.86	40.88
	0.8864	0.8844	0.8843	0.9858	0.9857	0.9990	0.7381	0.8769	0.9556	0.9871
IPAD	35.86	35.78	35.76	47.66	47.62	58.62	29.86	35.74	41.18	46.43
	0.9457	0.9452	0.9451	0.9903	0.9902	0.9989	0.8624	0.9378	0.9743	0.9903
BitNet (0.94M)	39.34	39.49	39.49	49.72	49.68	57.55	33.46	39.12	44.02	48.46
	0.9701	0.9719	0.9719	0.9954	0.9954	0.9989	0.9128	0.9632	0.9853	0.9943
BE-CALF (5.18M)	39.91	39.98	39.98	51.14	51.14	59.51	-	-	-	-
	0.9737	0.9752	0.9752	0.9940	0.9940	0.9993	-	-	-	-
D16 (<15.46M)	41.19	41.51	41.51	53.47	53.48	63.51	34.11	39.95	44.94	49.72
	0.9794	0.9810	0.9810	0.9980	0.9979	0.9998	0.9279	0.9693	0.9876	0.9953
<i>RDN-ABCD (Ours)</i> (11.52M)	42.31	42.84	42.84	54.07	54.10	63.75	35.14	40.94	45.68	50.08
	0.9831	0.9847	0.9847	0.9984	0.9984	0.9998	0.9392	0.9746	0.9893	0.9957
<i>EDSR-ABCD (Ours)</i> (12.22M)	42.47	43.02	43.02	54.15	54.18	63.78	35.25	41.04	45.74	50.11
	0.9837	0.9852	0.9852	0.9984	0.9984	0.9998	0.9401	0.9748	0.9893	0.9957
<i>SwinIR-ABCD (Ours)</i> (12.10M)	42.51	43.03	43.03	54.08	54.12	63.74	35.44	41.18	45.80	50.13
	0.9844	0.9855	0.9855	0.9984	0.9984	0.9998	0.9412	0.9751	0.9895	0.9957

# Qualitative Comparison

(3-bit→8-bit)

**IPAD**

**BitNet**

**D16**

**ABCD(Ours)**

**GT**



# Qualitative Comparison

(3-bit→8-bit)

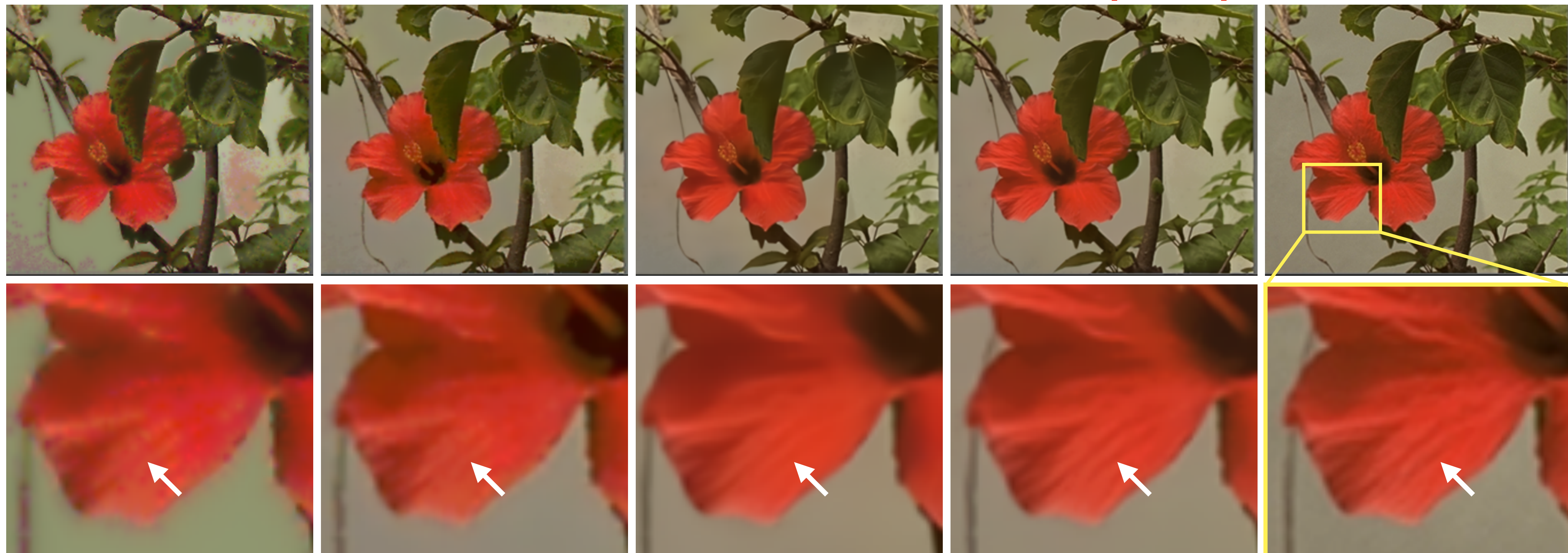
**IPAD**

**BitNet**

**D16**

**ABCD(Ours)**

**GT**



# Qualitative Comparison

(3-bit→8-bit)

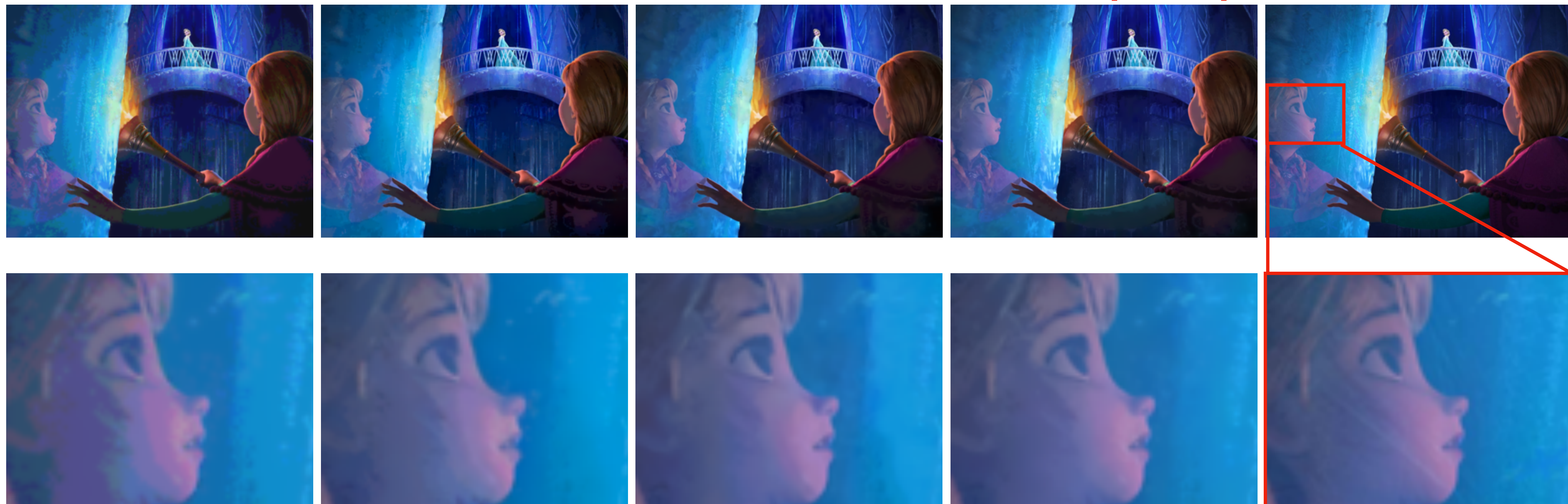
**IPAD**

**BitNet**

**D16**

**ABCD(Ours)**

**GT**



# Qualitative Comparison

(3-bit→8-bit)

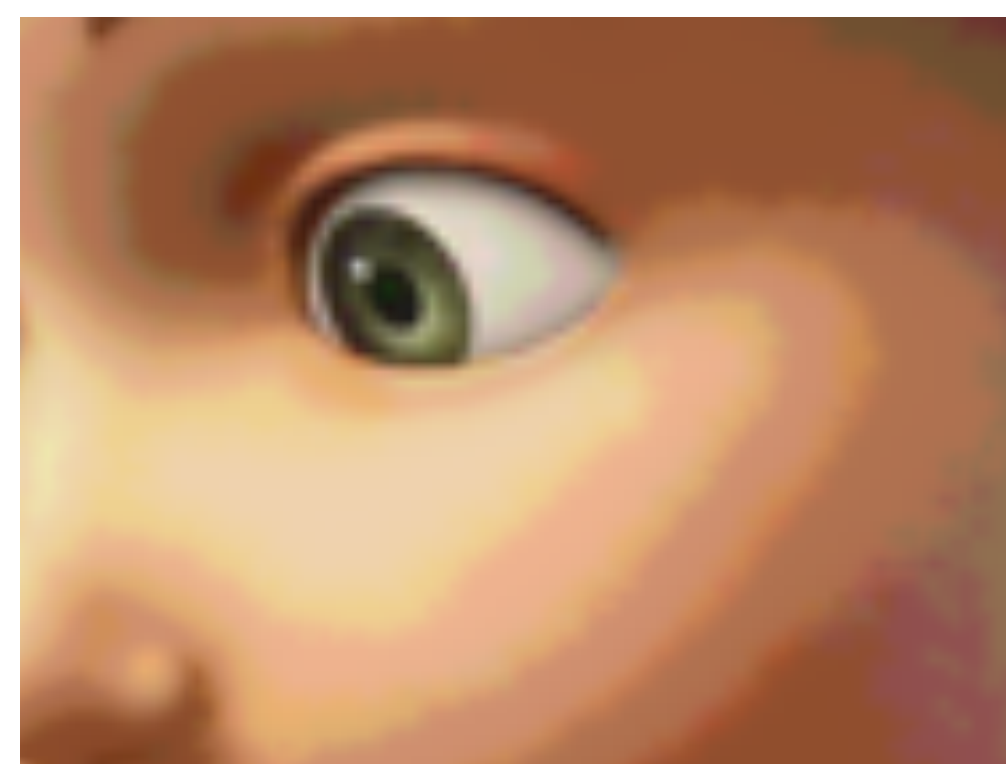
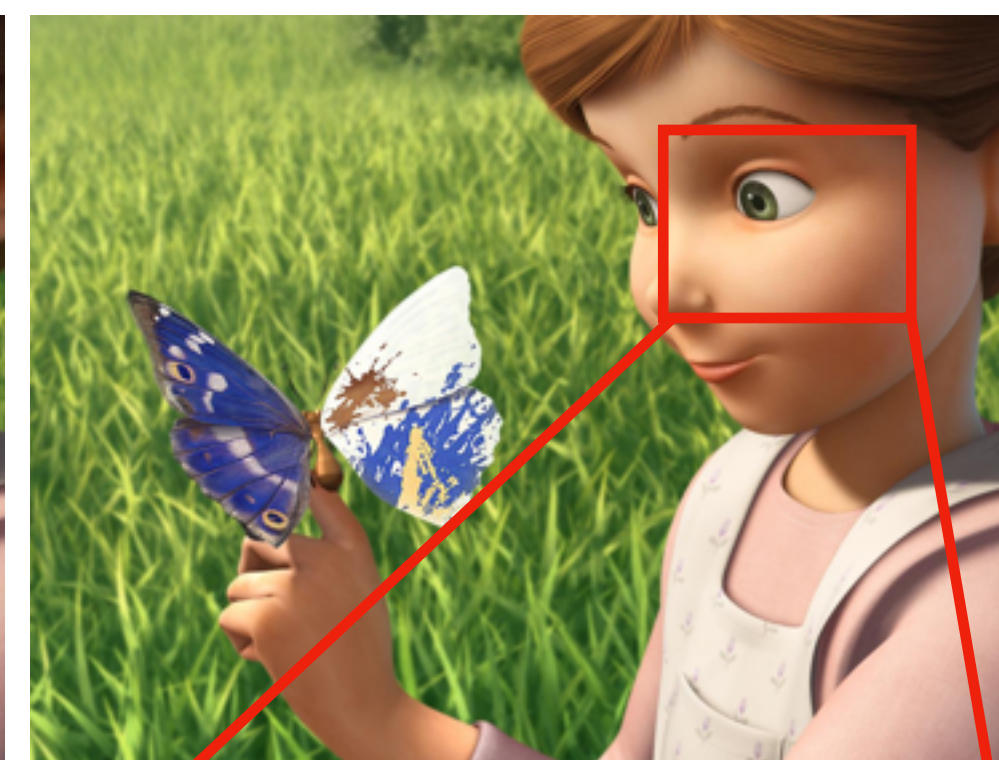
**IPAD**

**BitNet**

**D16**

**ABCD(Ours)**

**GT**



# Qualitative Comparison

(3-bit→8-bit)

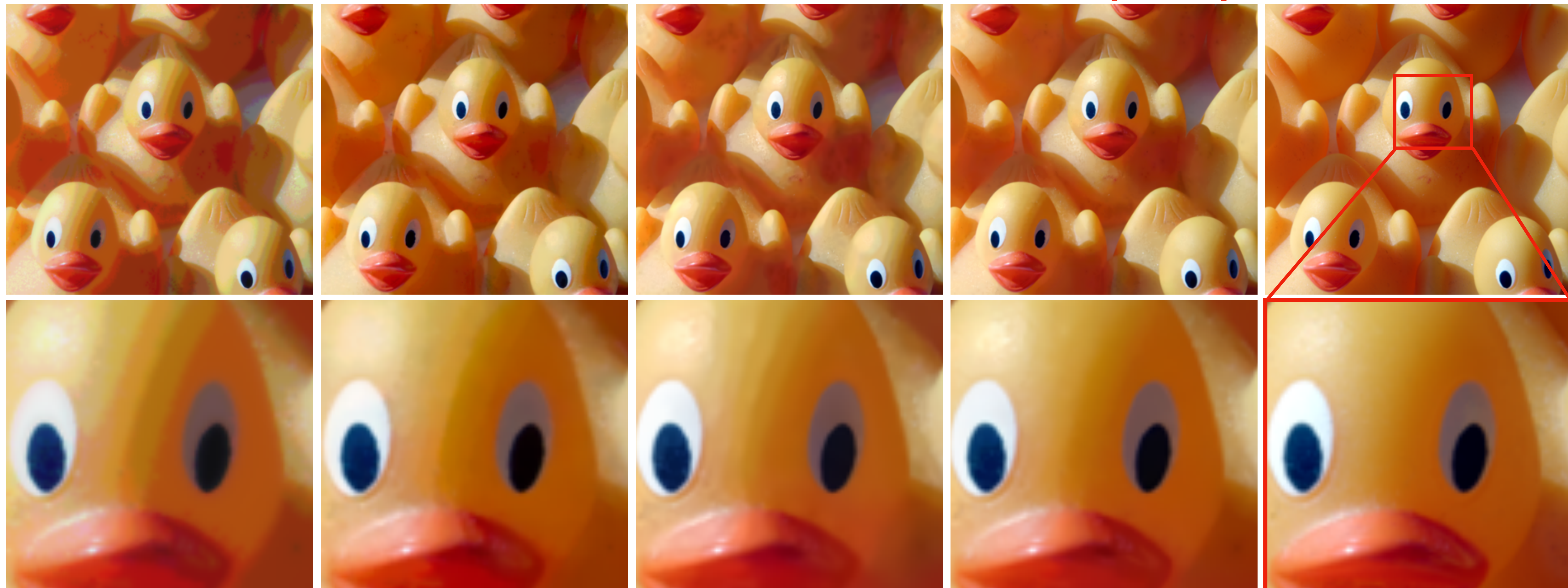
**IPAD**

**BitNet**

**D16**

**ABCD(Ours)**

**GT**





# Qualitative Comparison

(2-bit→8-bit) : Out-of-distribution

Input

IPAD

BitNet

**ABCD(Ours)**

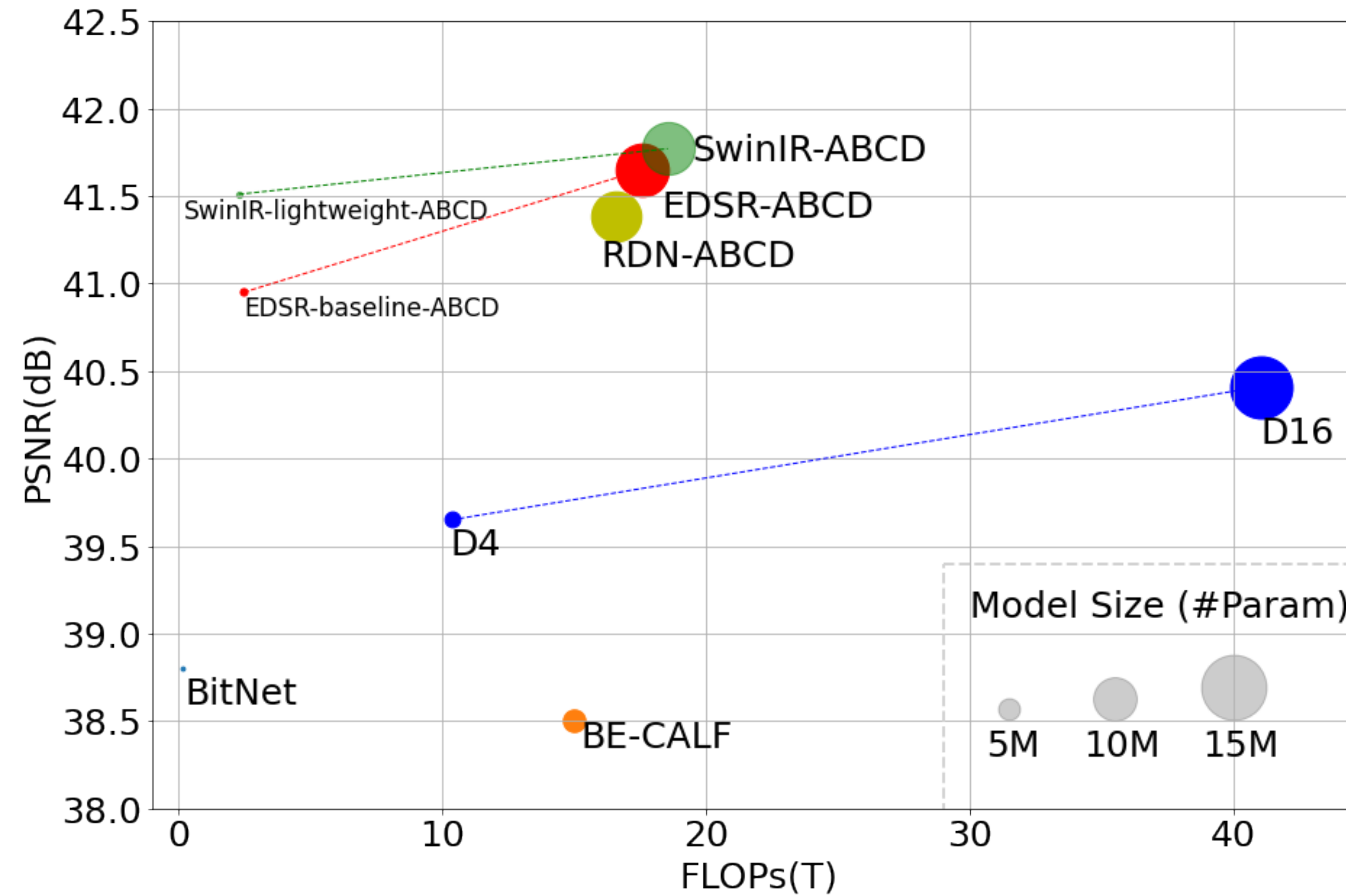
GT



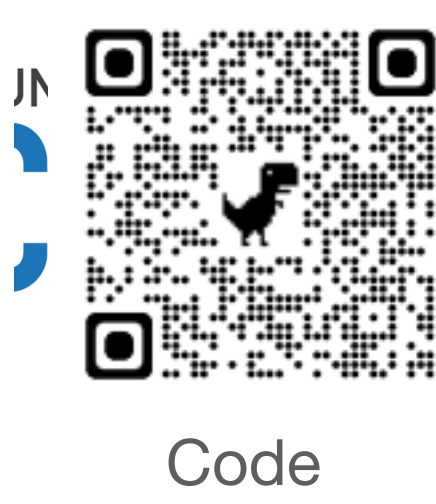
- 2-bit input is out-of-distribution cases for all methods

# FLOPs and Memories

## FLOPs vs PSNR



# TUE-PM-167



## 2-Bit

## 3-Bit

## 4-Bit

INPUT



ABCD



GT

