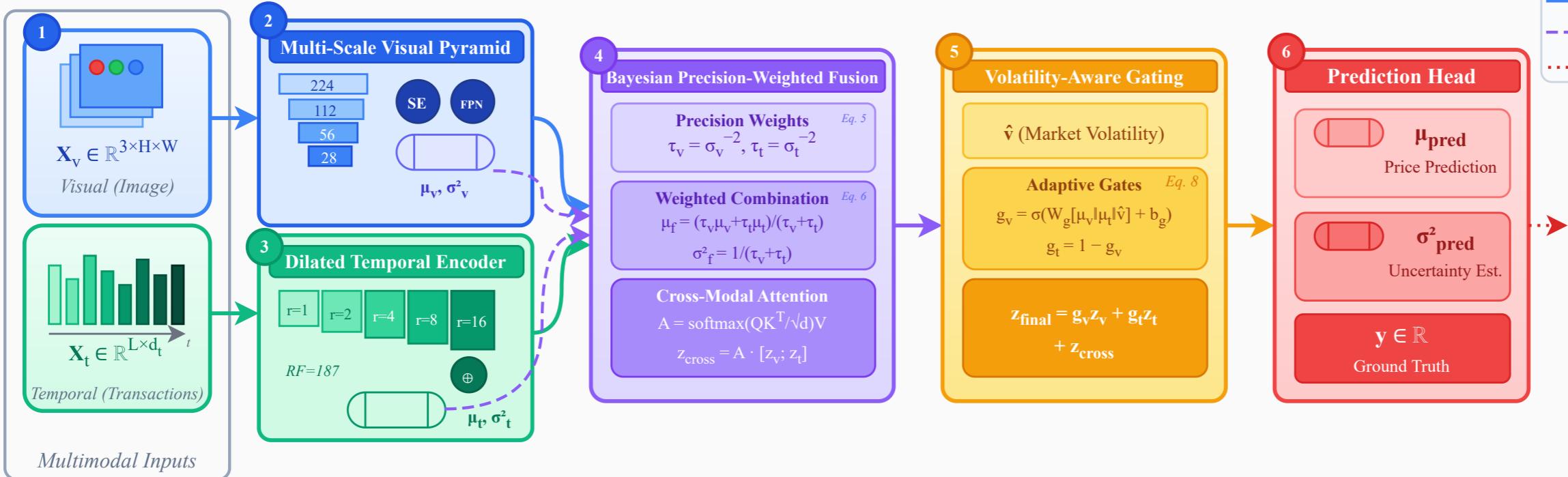


### (a) Complete UIBFuse Pipeline



### (b) Multi-Scale Visual Pyramid

Level 1: 224×224, C=64

Level 2: 112×112, C=128

Level 3: 56×56, C=256

Level 4: 28×28, C=512

#### SE-Net Attention

Squeeze: GAP → FC(C/r)  
Excitation: FC(C) →  $\sigma$   
Scale:  $x \cdot \sigma(s)$

#### FPN Multi-Scale Fusion

Top-down pathway  
Lateral connections  
Element-wise addition

**Output:**  $\mu_v \in \mathbb{R}^{256}, \sigma^2_v \in \mathbb{R}^{256}_{>0}$

*Spectral:  $S(f) \propto f^{-2.1} \rightarrow 4 \text{ levels capture } 95\% \text{ info}$*

### (c) Dilated Temporal Encoder

$X_t \in \mathbb{R}^{L \times 64}$

RF = 1 +  $\sum_i (k-1) \cdot r_i = 187 \text{ steps}$

#### Dilated Convolution Stack (Octave-based Rates)

r=1  
( $2^0$ )

r=2  
( $2^1$ )

r=4  
( $2^2$ )

r=8  
( $2^3$ )

r=16  
( $2^4$ )

Transformer Layers (H=16 heads,  $d_{\text{model}}=256$ )

**Output:**  $\mu_t \in \mathbb{R}^{256}, \sigma^2_t \in \mathbb{R}^{256}_{>0}$

### (d) Information-Theoretic Derivations

#### Information Bottleneck Objective

$$\mathcal{L}_{IB} = -I(Z; Y) + \beta \cdot I(Z; X_v, X_t)$$

#### ① Latent Dimension ( $d_z$ )

$$d_z \geq 2 \cdot I((V, T); Y) / (\beta \cdot \log(2\pi e))$$

$I \approx 4.3 \text{ nats}, \beta = 0.1 \rightarrow d_z = 256$

#### ② Pyramid Scales

$$S(f) \propto f^{-\alpha}, \alpha \approx 2.1 \rightarrow 95\% \text{ info at } \{224, 112, 56, 28\}$$

#### ③ Attention Heads (H)

Spectral clustering → 10-20 modes →  $H = 16$