

Lossless Video Coding Based on Probability Model Optimization Utilizing Example Search and Adaptive Prediction

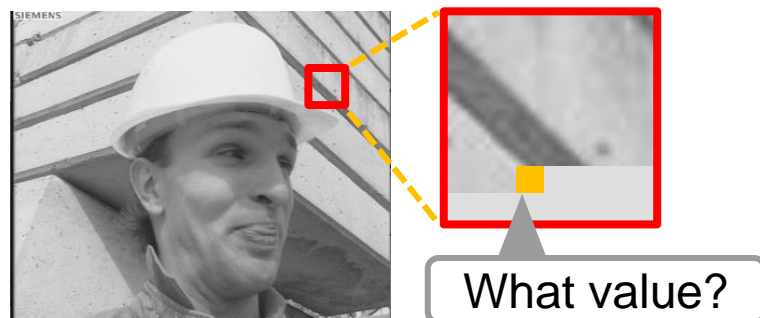
**Kyohei Unno^{†‡}, Koji Nemoto[‡], Yusuke Kameda[‡],
Ichiro Matsuda[‡], Susumu Itoh[‡], Sei Naito[†]**

† KDDI Research, Inc.

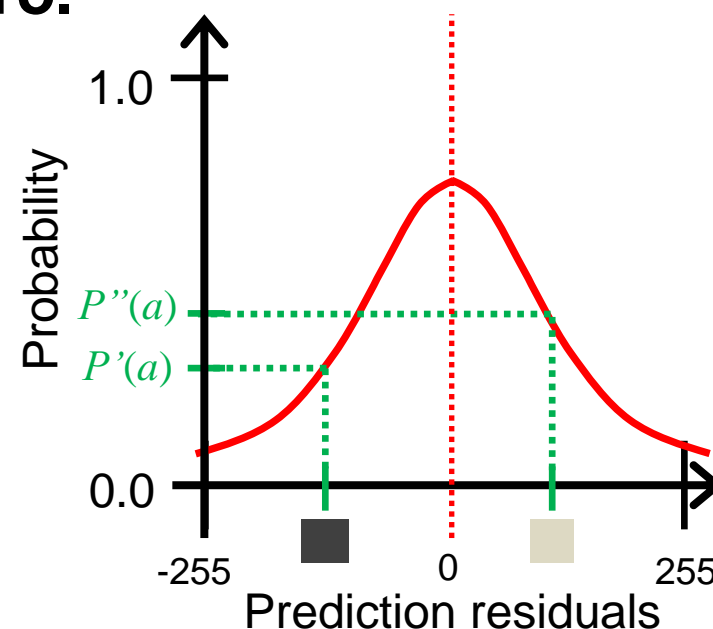
‡ Tokyo University of Science

Presenter : Kyohei Unno

- Conventional video coding methods typically consist of two stages.
 1. De-correlation (e.g. Prediction to get residuals)
 2. Entropy coding (e.g. Arithmetic coding)
- In general, probabilities of de-correlated signals are modeled as single-peaked symmetric functions centered at zero.



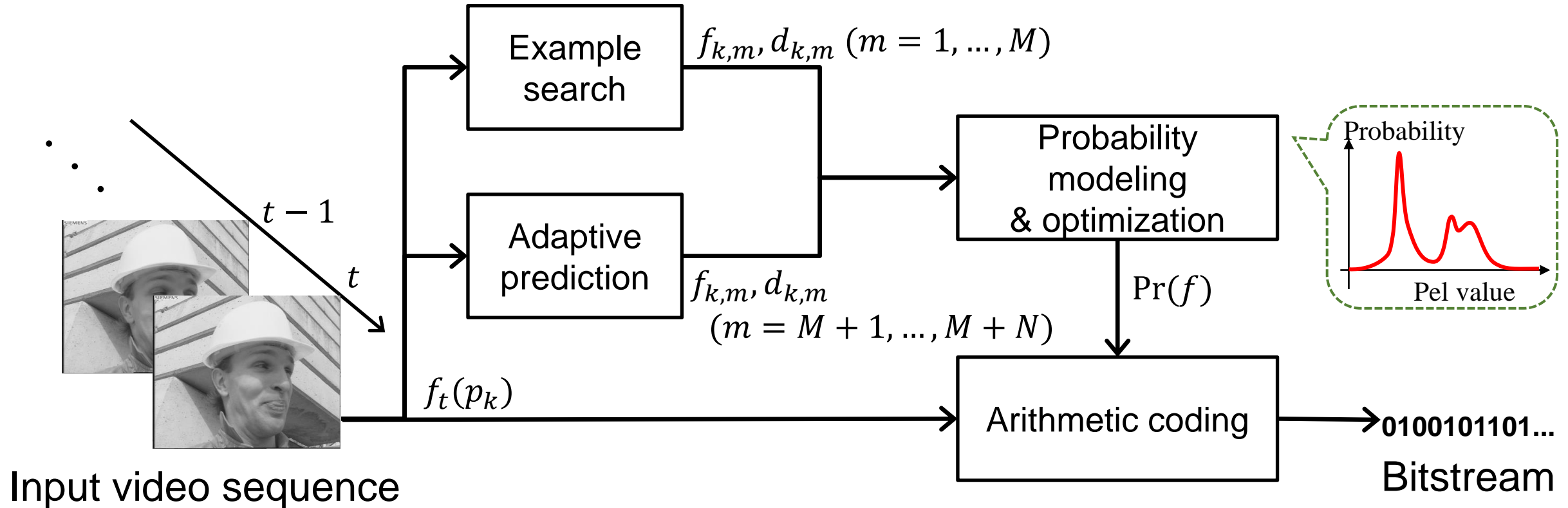
Actual value may be  or 



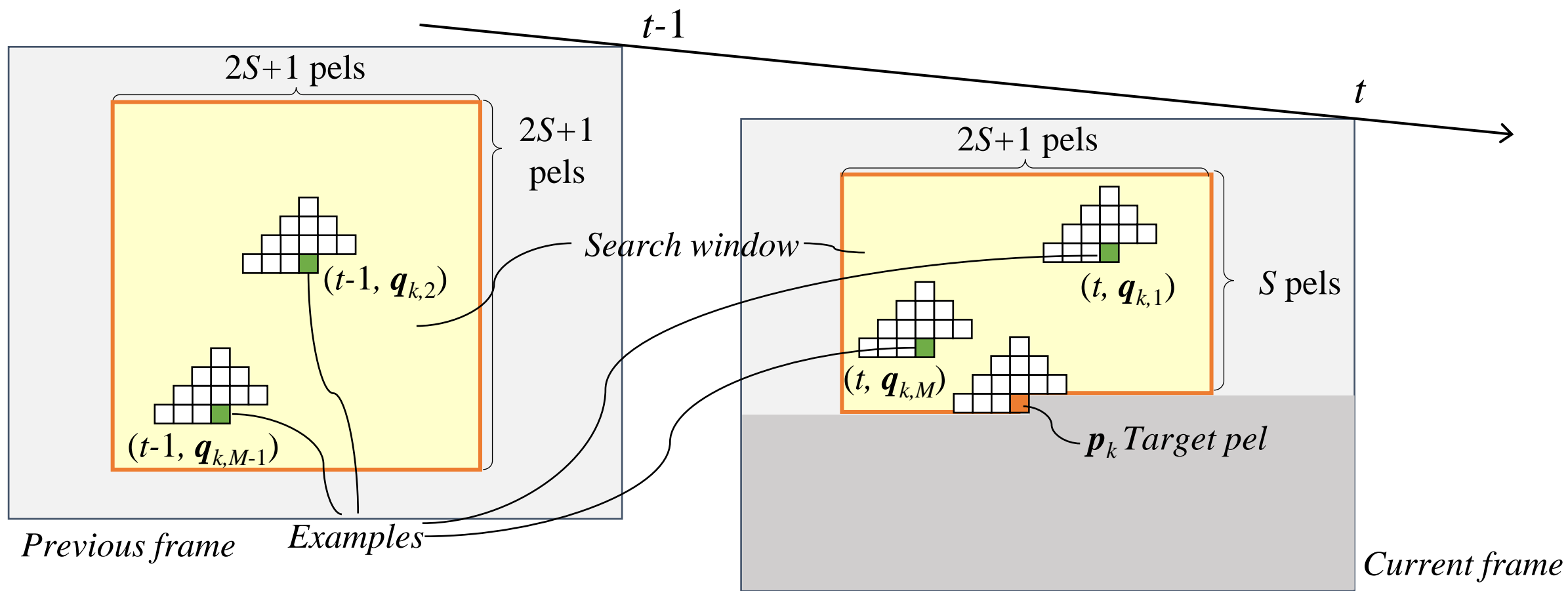
Such a model can not work well at texture boundaries 😞.

Overview of the proposed method

- Proposed method estimates probability of pel values directly with multi-peaked Gaussian mixture model (GMM) pel-by-pel.
- Center position $f_{k,m}$ and reliability $d_{k,m}$ (related to variance) of each Gaussian are estimated by **example search** and **adaptive prediction**.



- M pels $\{q_{k,1}, \dots, q_{k,M}\}$ that have smaller template matching costs are collected as examples from the search window.
- The search window is set to both the current frame and the previous frame.



- Center position $f_{k,m}$ and reliability metric $d_{k,m}$ for each gaussian are given by the examples and the corresponding cost, respectively

- Center position $f_{k,m}$: Example pel value with local mean compensation.

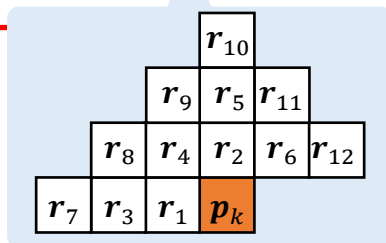
$$f_{k,m} = f_{t-\tau}(\mathbf{q}_{k,m}) - \mu_{t-\tau}(\mathbf{q}_{k,m}) + \mu_t(\mathbf{p}_k)$$

Pel value of $\mathbf{q}_{k,m}$

Local means of neighboring pels at \mathbf{p}_k and $\mathbf{q}_{k,m}$

- Reliability metric $d_{k,m}$: Template matching cost of the k -th example.

$$d_{k,m} = \left[\sum_{i=1}^{12} w_i \cdot \left(f_{t-\tau}(\mathbf{q}_{k,m} + \mathbf{r}_i) - \mu_{t-\tau}(\mathbf{q}_{k,m}) - f_t(\mathbf{p}_k + \mathbf{r}_i) + \mu_t(\mathbf{p}_k) \right)^2 \right]^{\frac{1}{2}} + \lambda_d \cdot \|\mathbf{q}_{k,m} - \mathbf{p}_k\|_1$$

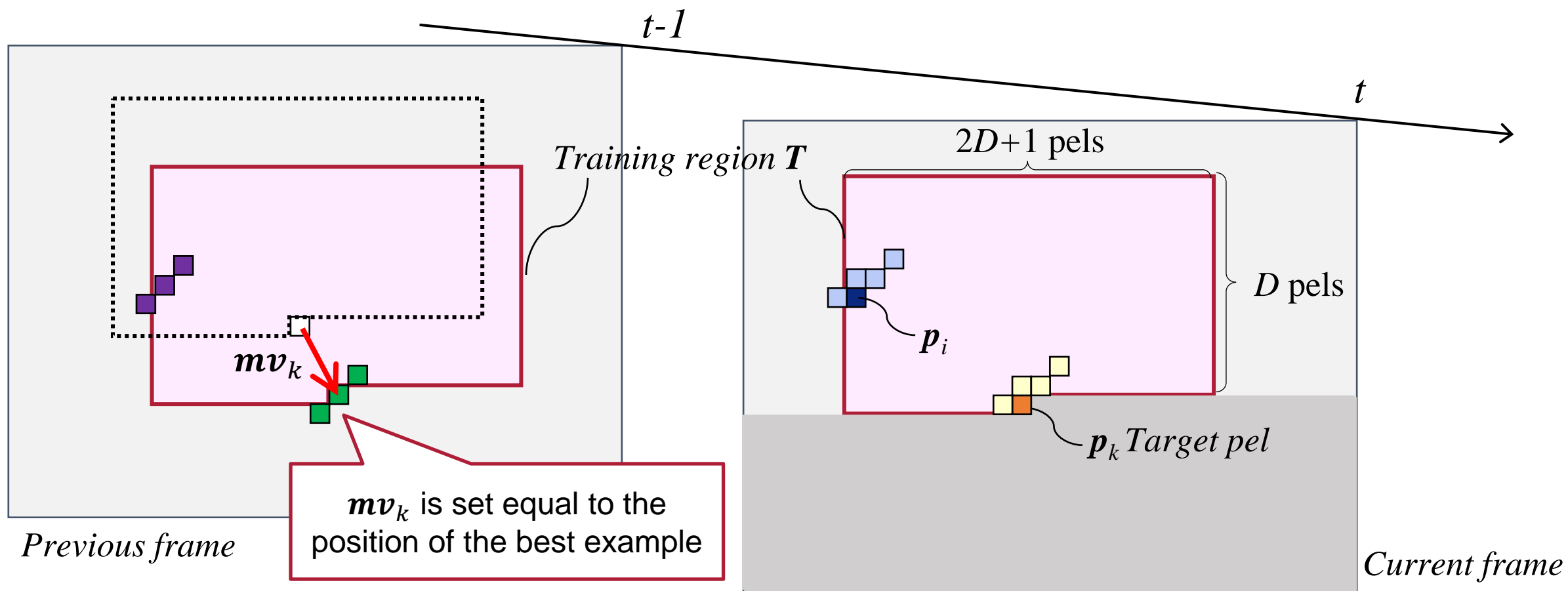


Weighted SSD with local mean compensation

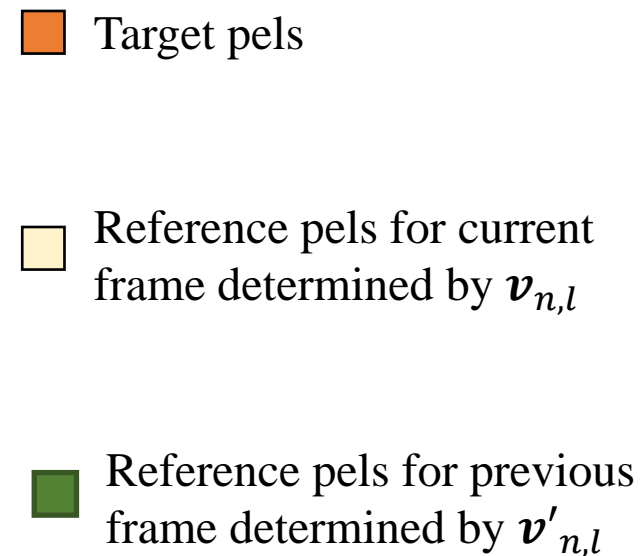
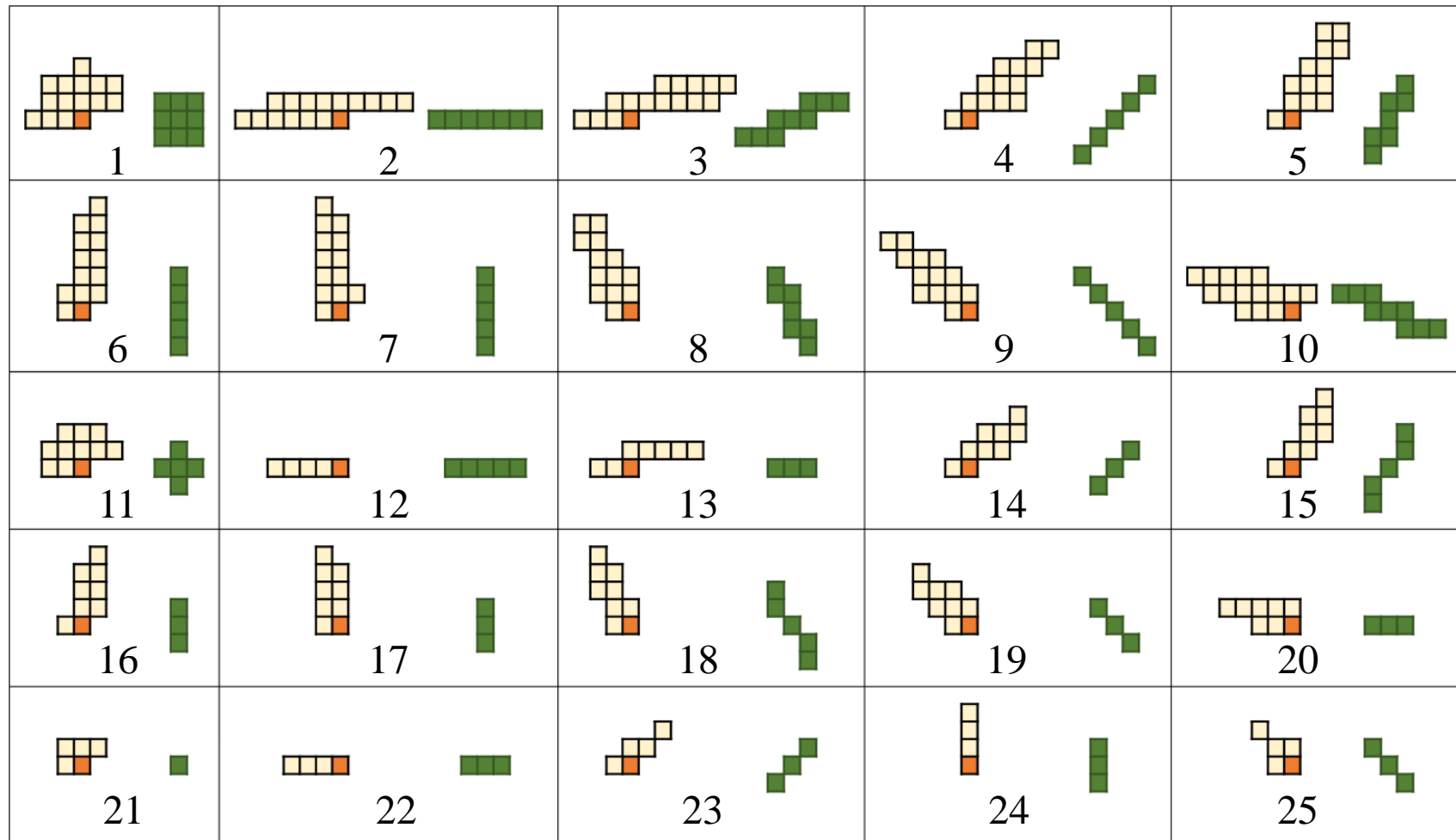
Regularization term

- The number of examples M is optimized for every region of 64 x 64 pels.

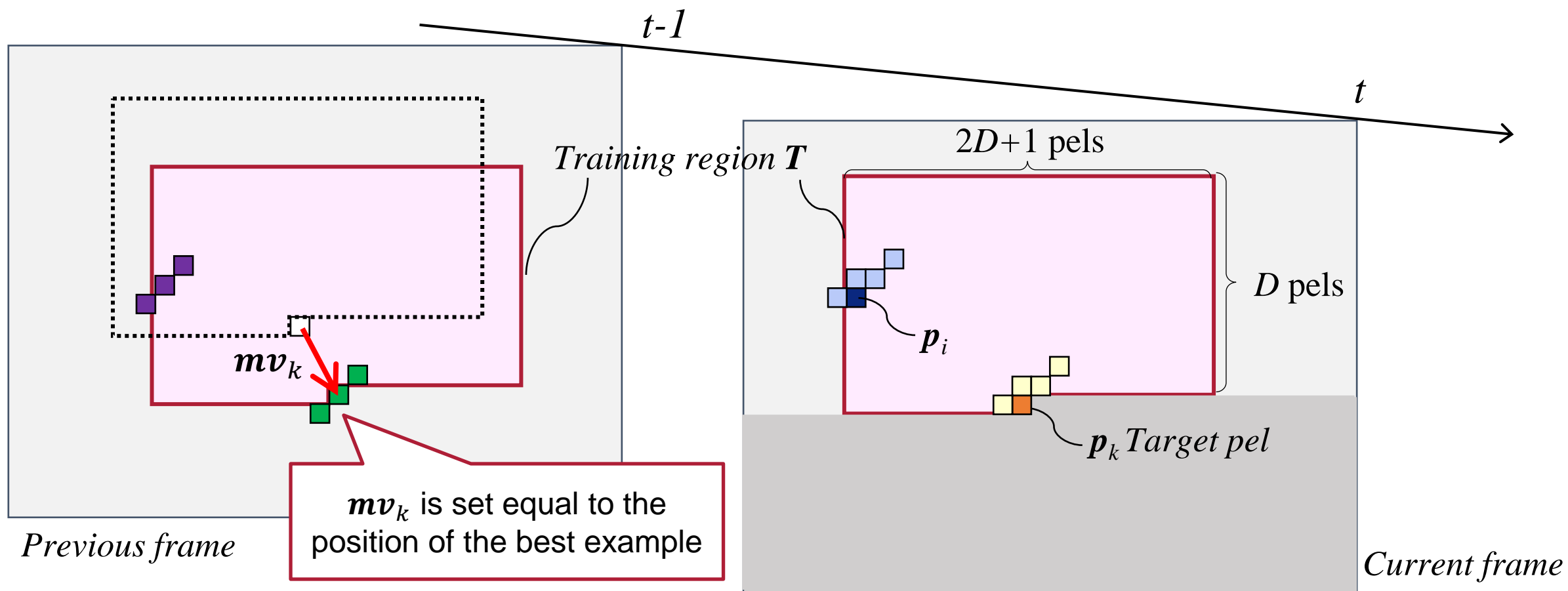
- N kinds of affine predictors (i.e. linear + constant terms) are trained pel-by-pel by the weighted least-square method.
 - Predictors refer pels in the current frame and the previous frame simultaneously.



- We use $N = 25$ kinds of predictors for each pel.
 - Each predictor has unique arrangement of reference pels.



- N kinds of affine predictors (i.e. linear + constant terms) are trained pel-by-pel by the weighted least-square method.
 - Predictors refer pels in the current frame and the previous frame simultaneously.



- Center position $f_{k,m}$ and reliability metric $d_{k,m}$ for each gaussian are given as the predicted value and the training error, respectively

- Center position $f_{k,m}$: Predicted value.

$$f_{k,M+n} = b_{n,0} + \sum_{l=1}^{L_n} b_{n,l} \cdot f_t(\mathbf{p}_k + \mathbf{v}_{n,l}) + \sum_{l=1}^{L'_n} b_{n,L_n+l} \cdot f_{t-1}(\mathbf{p}_k + m\mathbf{v}_k + \mathbf{v}'_{n,l})$$

Number of reference pels

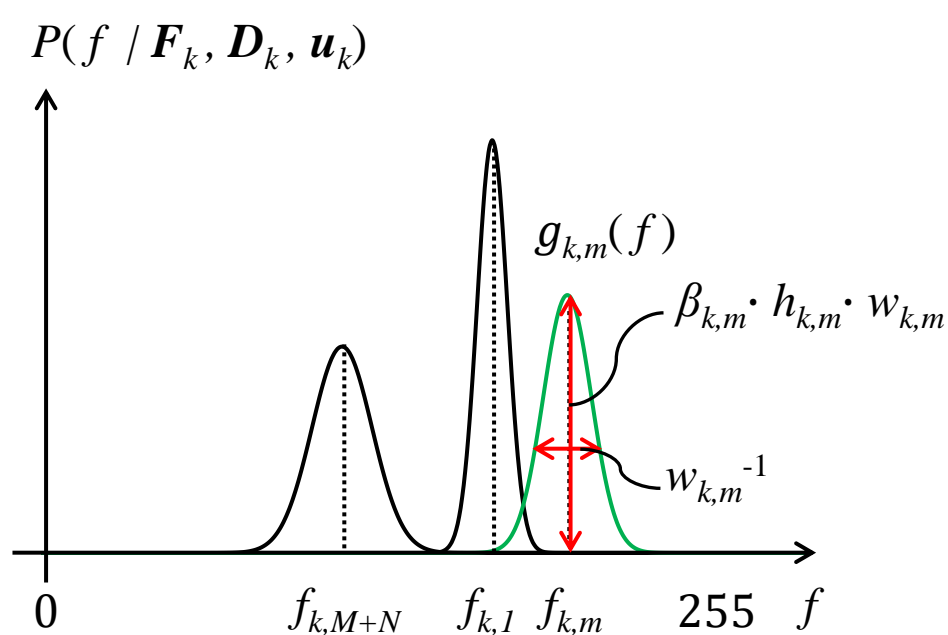
Arrangement of reference pels

- Reliability metric $d_{k,m}$: Training error.

$$d_{k,M+n} = \left(\sum_{\mathbf{p}_i \in T} \left((f_t(\mathbf{p}_i) - \hat{f}_t(\mathbf{p}_i, n))^2 / \sigma_{n,i}^2 \right) / \sum_{\mathbf{p}_i \in T} (1 / \sigma_{n,i}^2) \right)^{\frac{1}{2}}$$

Squared training prediction error by predictor n

Squared error between neighboring pels of \mathbf{p}_k and \mathbf{p}_i



$$P(f | \mathbf{F}_k, \mathbf{D}_k, \mathbf{u}_k) = \sum_m^{M+N} g_{k,m}(f) + \epsilon$$

$$g_{k,m} = \beta_{k,m} \cdot h_{k,m} \cdot w_{k,m} \cdot \exp(-w_{k,m}^2 \cdot (f - f_{k,m})),$$

$$h_{k,m} = \exp(-a_1 \cdot d_{k,m}),$$

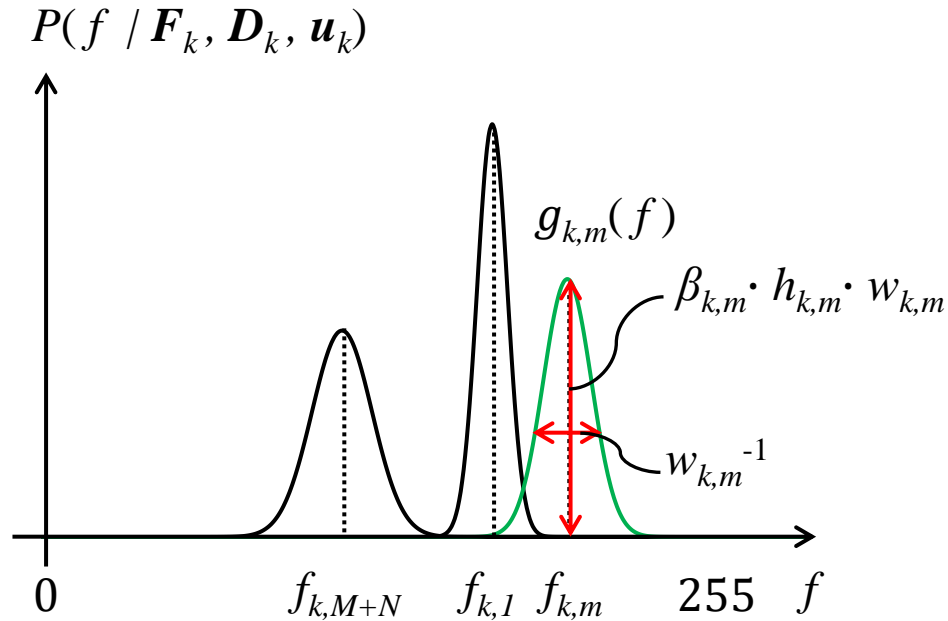
$$w_{k,m} = a_0 \cdot \exp(-a_2 \cdot d_{k,m}) \cdot \exp(-a_3 \cdot u_{k,m})$$

$$\text{Bitrate}(\mathbf{p}_k) = -\log_2 \frac{P(f(\mathbf{p}_k) | \mathbf{F}_k, \mathbf{D}_k, \mathbf{u}_k)}{\sum_{f=0}^{255} P(f | \mathbf{F}_k, \mathbf{D}_k, \mathbf{u}_k)}$$

■ **GMM is defined by center positions $\mathbf{F}_k = \{f_{k,m}\}$, reliability metrics $\mathbf{D}_k = \{d_{k,m}\}$, and local feature $\mathbf{u}_k = \{u_{k,m}\}$ ($m = 1, \dots, M + N$).**

- $\{f_{k,m}, d_{k,m}, u_{k,m} \mid m = 1, \dots, M\}$ come from example search.
- $\{f_{k,m}, d_{k,m}, u_{k,m} \mid m = M + 1, \dots, M + N\}$ come from adaptive prediction.
- \mathbf{u}_k is local average of bitrates in the causal area and used for context modeling.

■ **Shape of the probability model is controllable by model parameters a_0, \dots, a_3 .**



$$P(f | \mathbf{F}_k, \mathbf{D}_k, \mathbf{u}_k) = \sum_m^{M+N} g_{k,m}(f) + \epsilon$$

$$g_{k,m} = \beta_{k,m} \cdot h_{k,m} \cdot w_{k,m} \cdot \exp(-w_{k,m}^2 \cdot (f - f_{k,m})),$$

$$h_{k,m} = \exp(-a_1 \cdot d_{k,m}),$$

$$w_{k,m} = a_0 \cdot \exp(-a_2 \cdot d_{k,m}) \cdot \exp(-a_3 \cdot u_{k,m})$$

$$\text{Bitrate}(\mathbf{p}_k) = -\log_2 \frac{P(f(\mathbf{p}_k) | \mathbf{F}_k, \mathbf{D}_k, \mathbf{u}_k)}{\sum_{f=0}^{255} P(f | \mathbf{F}_k, \mathbf{D}_k, \mathbf{u}_k)}$$

- Model parameters a_0, \dots, a_3 and the number of examples M are optimized to minimize bitrate in every region of 64 x 64 pels.
 - a_0, \dots, a_3 are optimized by non-linear optimization technique (the quasi-Newton method).
 - Optimal setting of $M (= 0 \sim 31)$ is chosen in the optimization process.
 - The number of predictors $N (= 25)$ is fixed. Optimization of N is a future task.
- a_0, \dots, a_3 and M are sent to a decoder by fixed-length coding.

■ Test sequences

- **JVET CTC Class B, C, D, and E sequences.**
 - Align with Low-Delay P condition.
- **Converted from 10 bits to 8 bits by 2 bits right shift.**
- **Only luma signal is tested.**
- **The first 15 and 30 frames are used for Class B and the other classes, respectively.**
 - Due to the long processing time of the proposed method.

■ Comparison with lossless mode of HEVC RExt and VVC

- **HEVC RExt : HM 16.20**
- **VVC : VTM 8.2**
- **The number of reference pictures is restricted to 1 for both HM and VTM.**

JVET Test sequences

Class B
1920x1080



MarketPlace



RitualDance



Cactus



BasketBallDrive



BQTerrace

Class C
832x480



BasketBallDrill



BQMall



PartyScene



RaceHorses

Class D
416x240



BasketBallPass



BQSquare



BlowingBubbles



RaceHorses

Class E
1280x720



FourPeople



Johnny



KristenAndSara

■ Test sequences

- **JVET CTC Class B, C, D, and E sequences.**
 - Align with Low-Delay P condition.
- **Converted from 10 bits to 8 bits by 2 bits right shift.**
- **Only luma signal is tested.**
- **The first 15 and 30 frames are used for Class B and the other classes, respectively.**
 - Due to the long processing time of the proposed method.

■ Comparison with lossless mode of HEVC RExt and VVC

- **HEVC RExt : HM 16.20**
- **VVC : VTM 8.2**
- **The number of reference pictures is restricted to 1 for both HM and VTM.**

Comparison of coding rates (bits/pel).

■ Results

- Proposed method achieves about 9% and 5% better coding rates than HEVC RExt and VVC, respectively.
- Over 13% bitrate reductions from HEVC RExt are observed for RaceHorses (Class C) and all of Class E sequences.

■ Weakness of the proposed method is its processing time.

- Our implementation has not yet been tuned for the processing time.

| Test sequences | | Proposed | HEVC RExt | VVC |
|----------------------------------|------------------------|----------|-----------|--------|
| MarketPlace | Class B (1920x1080) | 3.229 | 3.675 | 3.532 |
| RitualDance | | 1.612 | 2.036 | 1.940 |
| Cactus | | 3.814 | 4.120 | 3.946 |
| BasketballDrive | | 3.473 | 3.785 | 3.602 |
| BQTerrace | | 3.647 | 3.879 | 3.749 |
| BasketballDrill | Class C (832x480) | 2.866 | 3.157 | 3.039 |
| BQMall | | 3.404 | 3.660 | 3.505 |
| PartyScene | | 3.828 | 3.973 | 3.794 |
| RaceHorses | | 3.329 | 3.902 | 3.793 |
| BasketballPass | Class D (416x240) | 2.068 | 2.182 | 2.114 |
| BQSquare | | 3.459 | 3.703 | 3.426 |
| BlowingBubbles | | 3.844 | 3.889 | 3.767 |
| RaceHorses | | 3.329 | 3.658 | 3.571 |
| FourPeople | Class E (1280x720) | 2.508 | 2.902 | 2.768 |
| Johnny | | 2.500 | 2.882 | 2.771 |
| KristenAndSara | | 2.376 | 2.741 | 2.618 |
| Average | | 3.080 | 3.384 | 3.246 |
| Bitrate reduction (vs HEVC RExt) | | -8.97% | 0.00% | -4.07% |

■ Proposed method

- Proposed method estimates probability of pel values directly with multi-peaked Gaussian mixture model (GMM) pel-by-pel.
- Center position $f_{k,m}$ and reliability $d_{k,m}$ (related to variance) of each Gaussian are estimated by example search and adaptive prediction.
- Model parameters of GMM a_0, \dots, a_3 and the number of examples M are optimized to minimize bitrate in every region of 64 x 64 pels.
- The proposed method achieves better bitrate in comparison with VVC and HEVC lossless mode.

■ Future works

- Reducing the processing time.
 - Adaptive prediction and optimization waste the most part of encoding time.
- Adaptive setting of the number of predictor N .
 - It might be effective for both improving coding performance and reducing processing time.

This work was supported by Ministry of Internal Affairs and Communications (MIC) of Japan (Grant no. JPJ000595).