

# Algorithm to capture local RNA sequential- and structural- motifs

Hiroshi Miyake

PhD student



THE UNIVERSITY OF TOKYO

Computational Approaches to RNA Structure and Function  
2018.7.15 – 7.27 @ Benasque

# Algorithm to capture local RNA sequential- and structural- motifs

Pico de Aneto

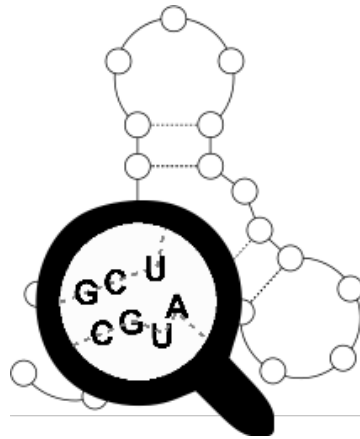


Great landscape !



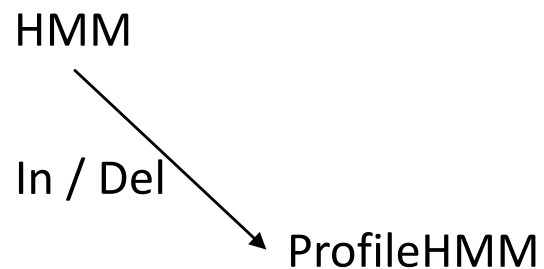
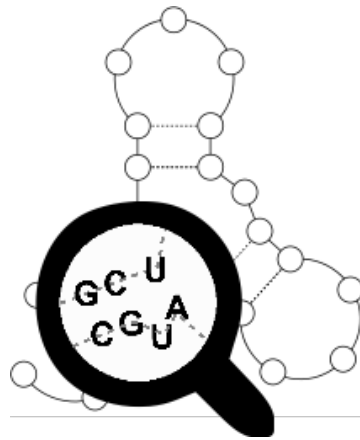
# RNA motif discovery

RNA motif = combination of conserved **subsequence** and **local structure**



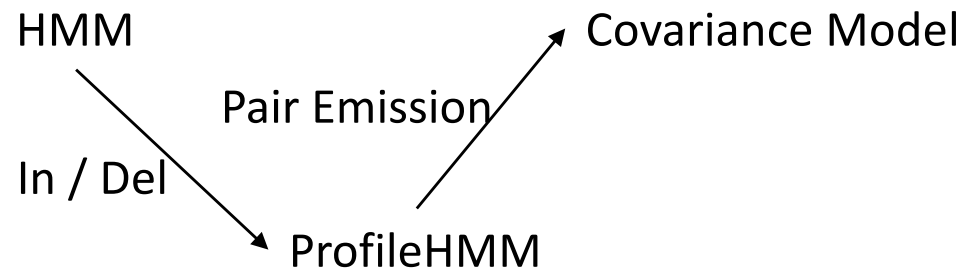
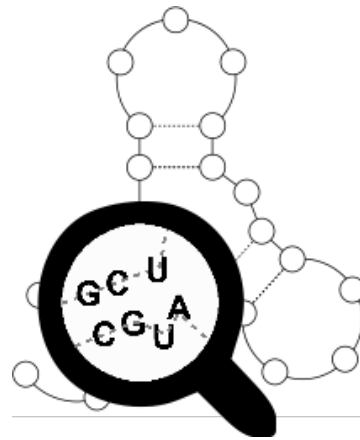
# RNA motif discovery

RNA motif = combination of conserved **subsequence** and **local structure**



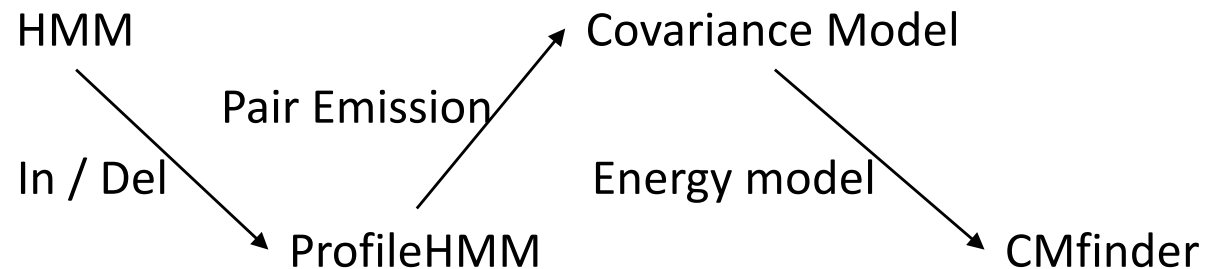
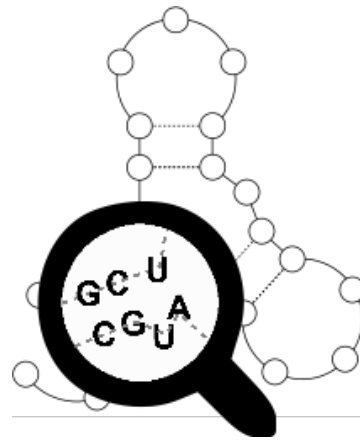
# RNA motif discovery

RNA motif = combination of conserved **subsequence** and **local structure**



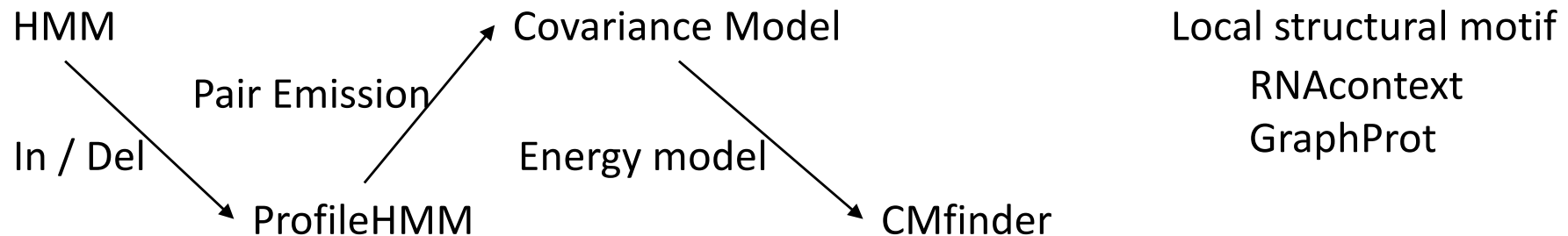
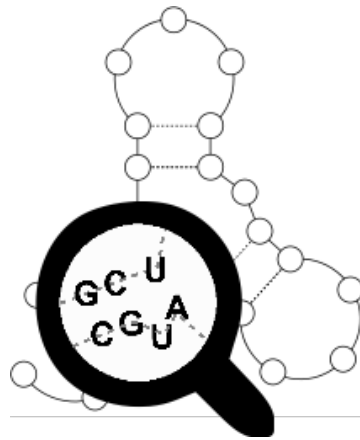
# RNA motif discovery

RNA motif = combination of conserved **subsequence** and **local structure**



# RNA motif discovery

RNA motif = combination of conserved **subsequence** and **local structure**



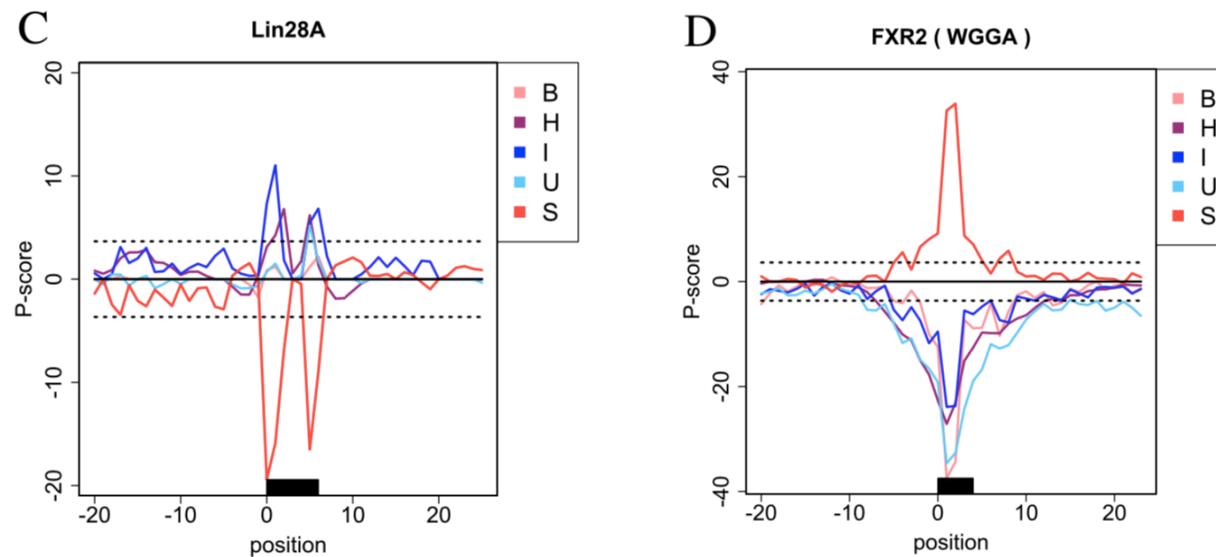
Kazan, H., Ray, D., Chan, E. T., Hughes, T. R., & Morris, Q. (2010). RNAcontext: a new method for learning the sequence and structure binding preferences of RNA-binding proteins. *PLoS computational biology*, 6(7), e1000832.

Maticzka, D., Lange, S. J., Costa, F., & Backofen, R. (2014). GraphProt: modeling binding preferences of RNA-binding proteins. *Genome biology*, 15(1), R17.

# Research aim

Open question: a mathematical model which can

- Capture complex 2D structural context



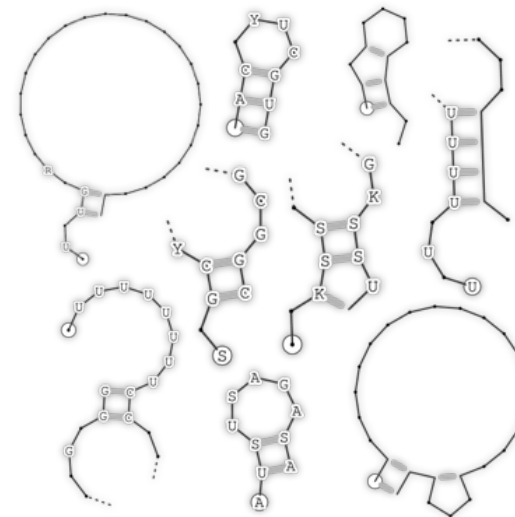
- Capture covariance between two loci
- Capture a gap inside motif



# Development of RNAelem iyak/RNAelem

input.fa

```
... GGGGAGGAAGTGGCTAGCTCAGGGCTTCAGGG...  
... ACAGACAGGGAGAGATGACTGAGTTAGATGAGA...  
... CGAGGGGGCGGGCTGGGGGTGCGAGAAGGAAGC...  
... TGGCAAGGAGACTAGGTCTAGGGGGACCACAGG...  
... GGGCAGGCTGCATGGAAAGGGGGCGGGCCTGG...  
... CTGCAGGCGGACCCCGTGAAGGGTTTCGCGGG...  
... TAGCGGGGACTCCTCGGGAGTCTTACAGGGCGG...  
... AGCTTAAGGTGCCGAAAAGTGGAAAATTACCA...  
... AAAGCAGGAAGGGAGGGTTAGCCTTGGGAAACC...  
... AATCTGGGTTTGCCACGGGGGCTTACTGAGTCA...  
... GGCCCCAGTCCCACAATTGGAAGAGATTGACG...  
... GTGTAGTGTCTTCAAGCTTGCTTTTTGGTGGGG...  
... ATTGGGGAGCTGTGGGGCGGCTGCCTTTGGTA...  
... GCTGTTGAGGGAGTCTGGGGCTTGTGAGCTGTA...
```



# SCFG for 2D structure

Discriminative model of 2D structure  $\sigma$  under given sequence  $x$

$$P(\sigma | x) = \frac{g(\sigma, x)}{Z(x)}, \quad Z(x) = \sum_{\sigma} g(\sigma, x)$$
$$g(\sigma, x) = e^{-\frac{1}{kT}\Delta G(\sigma, x)}$$

# SCFG for 2D structure

Discriminative model of 2D structure  $\sigma$  under given sequence  $x$

$$P(\sigma | x) = \frac{g(\sigma, x)}{Z(x)}, \quad Z(x) = \sum_{\sigma} g(\sigma, x)$$

$$g(\sigma, x) = e^{-\frac{1}{kT}\Delta G(\sigma, x)}$$

|              |   |  |
|--------------|---|--|
| $x$          | CAUGC <u>U</u> AGCUAGUCUGCGUAGCUGCGUACUAGCGCGUACGUCGGAU   |  |
| $\sigma \in$ | <ul style="list-style-type: none"> <li>.((( (. . (((((( . ((( . . . . . ))) . ))))))) . . . . .</li> <li>. . . (((((( . . . . . ))) . )))) . (((((( . . . . . ))) . . . . .</li> <li>. . . . . (((((( . ((( . ((( . . . . . ))) . ))))))) . ))))</li> <li>. . . . . (((((( . ((( . ((( . . . . . ))) . ))))))) . ))))</li> <li style="text-align: center;">⋮</li> </ul> | <ul style="list-style-type: none"> <li>-13.8</li> <li>-13.4</li> <li>-13.1</li> <li>-12.8</li> </ul> |

# SCFG for primary sequential profile

Discriminative model of alignment  $\psi$  with base- and base pair- emission, under given sequence  $x$  and pattern of interest

$$P(\psi | x) = \frac{h(\psi, x, \theta)}{Z(x, \theta)}, \quad Z(x, \theta) = \sum_{\psi} h(\psi, x, \theta)$$

# SCFG for primary sequential profile

Discriminative model of alignment  $\psi$  with base- and base pair- emission, under given sequence  $x$  and pattern of interest

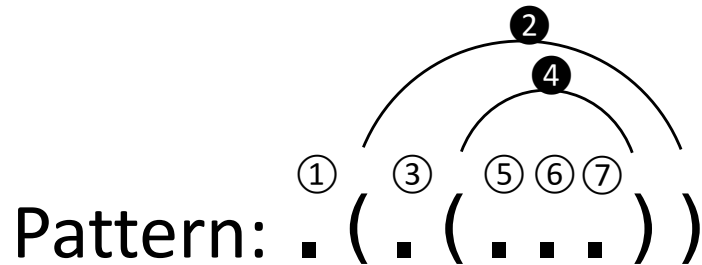
$$P(\psi | x) = \frac{h(\psi, x, \theta)}{Z(x, \theta)}, \quad Z(x, \theta) = \sum_{\psi} h(\psi, x, \theta)$$

$x$  CAUGCUAGCUAGUCUGCGUAGCUGCGUACUAGCGCGUACGUCGGAU

Pattern: . ( . ( . . . ) )

|            |                               | $h(\psi, x, \theta)$ |
|------------|-------------------------------|----------------------|
| $\psi \in$ | ***** . ( . ( . . . ) ) ***** | 0.6                  |
|            | ***** . ( . ( . . . ) ) ***** | 0.5                  |
|            | ***** . ( . ( . . . ) ) ***** | 0.5                  |
|            | ***** . ( . ( . . . ) ) ***** | 0.4                  |
|            | ⋮                             |                      |

# SCFG for primary sequential profile



|   |               |   |   |   |
|---|---------------|---|---|---|
|   | A             | C | G | U |
| ① |               |   |   |   |
| ③ |               |   |   |   |
| ⑤ | Base emission |   |   |   |
| ⑥ |               |   |   |   |
| ⑦ |               |   |   |   |
| * | background    |   |   |   |

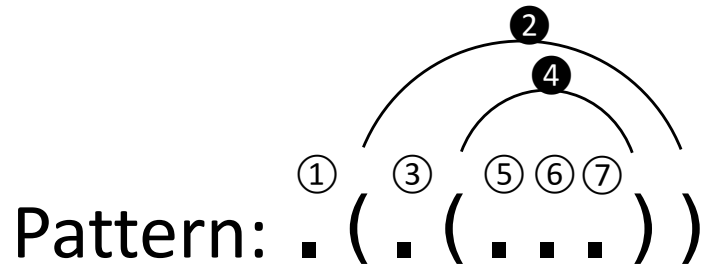
|   |                    |    |    |    |    |    |
|---|--------------------|----|----|----|----|----|
|   | CG                 | GC | AU | UA | GU | UG |
| ② | Base pair emission |    |    |    |    |    |
| ④ |                    |    |    |    |    |    |

$x$  CAUGCUAGCUAGUCUGCGUAGCUGCGUACUAGCGCGUACGUCGGAU

Pattern:  $\cdot ( \cdot ( \cdot \cdot \cdot ) )$

|            |   |   |     |
|------------|---|---|-----|
|            |   | $h(\psi, x, \theta)$                                |     |
| $\psi \in$ | { | ***** $\cdot ( \cdot ( \cdot \cdot \cdot ) )$ ***** | 0.6 |
|            |   | ***** $\cdot ( \cdot ( \cdot \cdot \cdot ) )$ ***** | 0.5 |
|            |   | ***** $\cdot ( \cdot ( \cdot \cdot \cdot ) )$ ***** | 0.5 |
|            |   | ***** $\cdot ( \cdot ( \cdot \cdot \cdot ) )$ ***** | 0.4 |
|            | ⋮ |   |     |

# SCFG for primary sequential profile



|   |               |   |   |   |
|---|---------------|---|---|---|
|   | A             | C | G | U |
| ① |               |   |   |   |
| ③ |               |   |   |   |
| ⑤ | Base emission |   |   |   |
| ⑥ |               |   |   |   |
| ⑦ |               |   |   |   |
| * | background    |   |   |   |

|   |                    |    |    |    |    |    |
|---|--------------------|----|----|----|----|----|
|   | CG                 | GC | AU | UA | GU | UG |
| ② | Base pair emission |    |    |    |    |    |
| ④ |                    |    |    |    |    |    |

Parameter  $\theta$

$x$  CAUGCUAGCUAGUCUGCGUAGCUGCGUACUAGCGCGUACGUCGGAU

Pattern:  $\cdot (\cdot (\cdot \cdot \cdot))$

|            |   |     |
|------------|---|-----|
| $\psi \in$ | ***** $\cdot (\cdot (\cdot \cdot \cdot))$ ***** | 0.6 |
|            | ***** $\cdot (\cdot (\cdot \cdot \cdot))$ ***** | 0.5 |
|            | ***** $\cdot (\cdot (\cdot \cdot \cdot))$ ***** | 0.5 |
|            | ***** $\cdot (\cdot (\cdot \cdot \cdot))$ ***** | 0.4 |

⋮

$h(\psi, x, \theta)$

# SCFG for primary sequential profile

Pattern:  $\cdot ( \cdot ( \cdot * \cdot \cdot ) )$

↑  
gap



# SCFG for primary sequential profile

Pattern:  $\cdot (\cdot (\cdot * \cdot \cdot))$

$x$  CAUGCUAGCUAGUCUGCGUAGCUGCGUACUAGCGCGUACGUCGGAU

Pattern:  $\cdot (\cdot (\cdot * \cdot \cdot))$

| $\psi \in$ | $h(\psi, x, \theta)$                                    | $h(\psi, x, \theta)$ |
|------------|---|----------------------|
| [          | ***** $\cdot (\cdot (\cdot$ ***** $\cdot \cdot))$ ***** | 0.6                  |
| [          | ***** $\cdot (\cdot (\cdot$ *** $\cdot \cdot))$ *****   | 0.5                  |
| [          | ***** $\cdot (\cdot (\cdot$ ***** $\cdot \cdot))$ ***** | 0.5                  |
| [          | ***** $\cdot (\cdot (\cdot$ ** $\cdot \cdot))$ ***      | 0.4                  |
|            | ⋮   |                      |

# Combined model

Joint probability of  $\psi$  and  $\sigma$ , under given sequence  $x$  and pattern of interest

$$P(\psi, \sigma | x) = \frac{f(\psi, \sigma, x, \theta, \lambda)}{Z(x, \theta, \lambda)}, \quad Z(x, \theta, \lambda) = \sum_{\psi} \sum_{\sigma} f(\psi, \sigma, x, \theta, \lambda)$$

$$f(\psi, \sigma, x, \theta, \lambda) = \mathbb{1}(\psi \circ \sigma) g(\sigma, x)^{\lambda} h(\psi, x, \theta)$$

$\mathbb{1}(\psi \circ \sigma)$  is 1 if paired loci are consistent among  $\psi$  and  $\sigma$ , and 0 otherwise.

$\lambda$  is a scalar, which can be also interpreted as “stability”

# Combined model

Joint probability of  $\psi$  and  $\sigma$ , under given sequence  $x$  and pattern of interest

$$P(\psi, \sigma | x) = \frac{f(\psi, \sigma, x, \theta, \lambda)}{Z(x, \theta, \lambda)}, \quad Z(x, \theta, \lambda) = \sum_{\psi} \sum_{\sigma} f(\psi, \sigma, x, \theta, \lambda)$$

$$f(\psi, \sigma, x, \theta, \lambda) = \mathbb{1}(\psi \circ \sigma) g(\sigma, x)^\lambda h(\psi, x, \theta)$$

$\mathbb{1}(\psi \circ \sigma)$  is 1 if paired loci are consistent among  $\psi$  and  $\sigma$ , and 0 otherwise.

$\lambda$  is a scalar, which can be also interpreted as “stability”

$$g(\sigma, x)^\lambda = \exp\left(-\frac{\lambda}{kT} \Delta G(\sigma, x)\right)$$

# Objective function and parameter fitting

Probability of motif existence

$$\begin{aligned} P(z = 0 \mid x; \theta, \lambda) &= \sum_{\sigma} P(\psi_0, \sigma \mid x; \theta, \lambda) \\ &= \sum_{\sigma} \frac{f(\psi_0, \sigma, x, \theta, \lambda)}{Z(x, \theta, \lambda)} \\ &= \sum_{\sigma} \frac{f(\psi_0, \sigma, x, \theta, \lambda)}{\sum_{\psi} \sum_{\sigma'} f(\psi, \sigma', x, \theta, \lambda)} \end{aligned}$$

$$P(z = 1 \mid x; \theta, \lambda) = 1 - P(z = 0 \mid x; \theta, \lambda)$$

$\psi_0$ : all bases are emitted by background state

$z \in \{0,1\}$ : motif existence in a sequence

# Objective function and parameter fitting

Probability of motif existence

$$\begin{aligned} P(z = 0 \mid x; \theta, \lambda) &= \sum P(\psi_0, \sigma \mid x; \theta, \lambda) \\ &= \sum_{\sigma} \frac{f(\psi_0, \sigma, x, \theta, \lambda)}{Z(x, \theta, \lambda)} \\ &= \sum_{\sigma} \frac{f(\psi_0, \sigma, x, \theta, \lambda)}{\sum_{\psi} \sum_{\sigma'} f(\psi, \sigma', x, \theta, \lambda)} \end{aligned}$$

$$P(z = 1 \mid x; \theta, \lambda) = 1 - P(z = 0 \mid x; \theta, \lambda)$$

$\psi_0$ : all bases are emitted by background state

$z \in \{0,1\}$ : motif existence in a sequence

Log likelihood over  
Positive / negative  
sequences

$$\begin{aligned} \mathcal{L}(\theta, \lambda) &= \sum_{x^+} \ln P(z = 1 \mid x^+; \theta, \lambda) + \sum_{x^-} \ln P(z = 0 \mid x^-; \theta, \lambda) \\ &= \mathcal{L}^+ + \mathcal{L}^- \end{aligned}$$


# Objective function and parameter fitting

Task:

$$\arg_{\theta, \lambda} \max \mathcal{L}(\theta, \lambda)$$

We can calculate exact derivation  $\left( \frac{\partial \mathcal{L}(\theta, \lambda)}{\partial \theta}, \frac{\partial \mathcal{L}(\theta, \lambda)}{\partial \lambda} \right)$  by nested inside-outside algorithm (ref. Sankoff's algorithm)

Expected value calculation  
DP over 2 SCFGs traversal !!


$$\frac{\partial L(\theta, \lambda)}{\partial \theta_{ij}} = \sum_{x^+} E[\mathcal{N}_{ij}(\psi)] \dots$$
$$\frac{\partial L(\theta, \lambda)}{\partial \lambda} = \sum_{x^-} E[\mathcal{G}(\sigma)] \dots$$

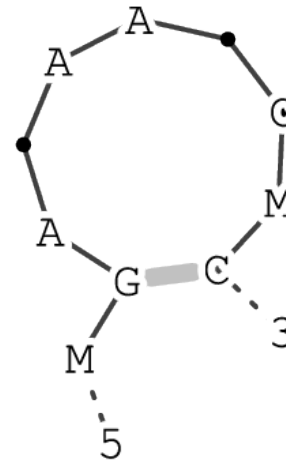
# Select the best pattern after optimization

1. Enumerate 2D structural patterns
2. Optimize parameter  $\theta$  and  $\lambda$  for each pattern
3. Model selection (k-fold cross validation) to select the best pattern

# Select the best pattern after optimization

1. Enumerate 2D structural patterns
2. Optimize parameter  $\theta$  and  $\lambda$  for each pattern
3. Model selection (k-fold cross validation) to select the best pattern

Output example





# Conclusion

- We formulated a combined motif model of  
primary sequence × secondary structure
- Traversing all the RNA 2D structural space enabled more precise prediction of the local structural motif
- The new model can assess the “local stability” of the 2D structure at binding region
- Several validated structural motifs were reproduced

# Acknowledgements



- The University Of Tokyo
  - Prof. Hisanori Kiryu
  - Prof. Kiyoshi Asai
  - Lab members
- AI Research Center (AIST)
  - Dr. Risa Kawaguchi

