

10 February 2025

# RNA Folding

## Plan.

- \* RNA Folding Problem
- \* Announcements
- \* DP for RNA Folding.

Another problem from biology

RNA Sequences  $m \in \{A, U, C, G\}^*$

$m = A U U C G G A C G G A A$

Key Question:

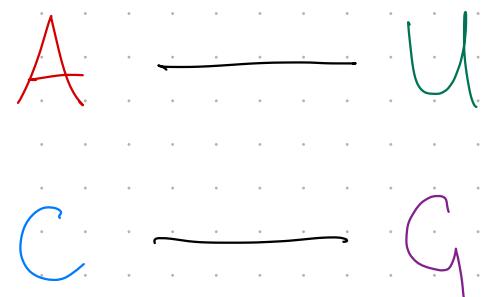
What structure does the RNA molecule adopt?

What structure does the RNA molecule adopt?

\_\_\_\_\_



Base pairing



} allowable pairings

What structure does the RNA molecule adopt?

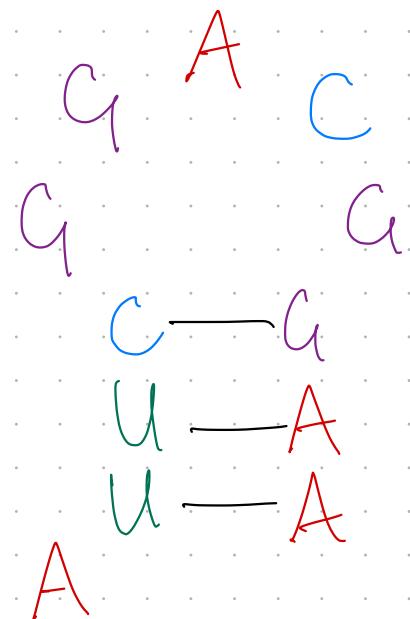
m = A U U C G G A C G G A A

↓ base pairing

What structure does the RNA molecule adopt?

m = A U U C G G A C G G A A

↓ base pairing



Given. RNA Sequence me  $\{A, U, C, G\}^*$

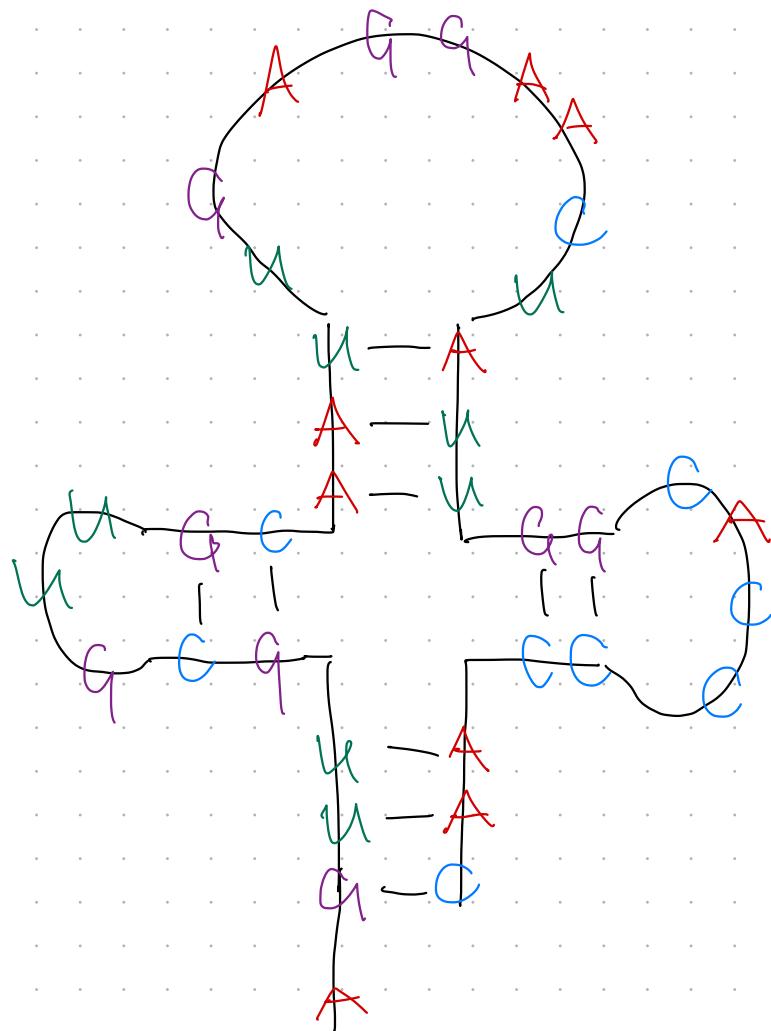
Find. Minimum "energy" structure

AquuAGCQuUGCAAUUGAGGAAACUAAUUGGCACCCCCAAC

Given. RNA Sequence me  $\{A, U, C, G\}^*$

Find. Minimum "energy" structure

A G U U A G C G U U G C A A U U G A G G A A C U A U U G G C A C C C C A A C



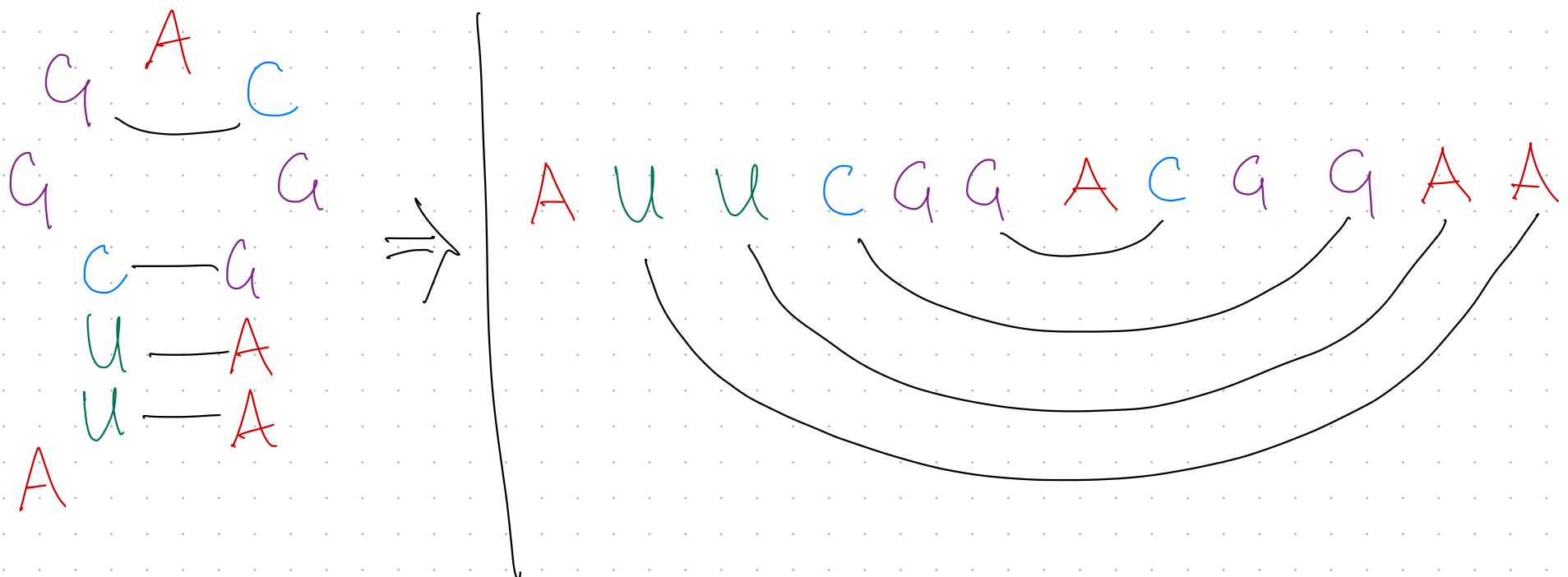
Min "energy"

$\equiv$

Max "legal" base pairing

Given. RNA Sequence me  $\{A, U, C, G\}^*$

Compute. Maximum Non-crossing matching E  
of preferred base pairs within m.

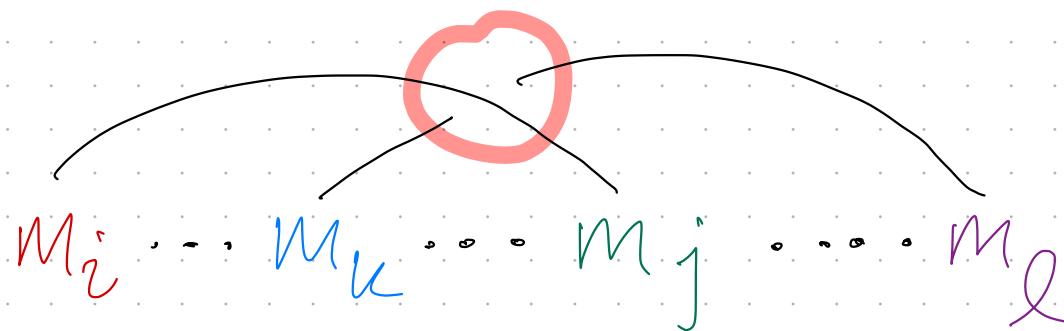


Given. RNA Sequence  $m \in \{A, U, C, G\}^*$

Compute. Maximum Non-crossing matching  $E$  of preferred base pairs within  $m$ .

Non-crossing. if  $(i, j), (k, l) \in E$

then,  $\neg (i < k < j < l)$



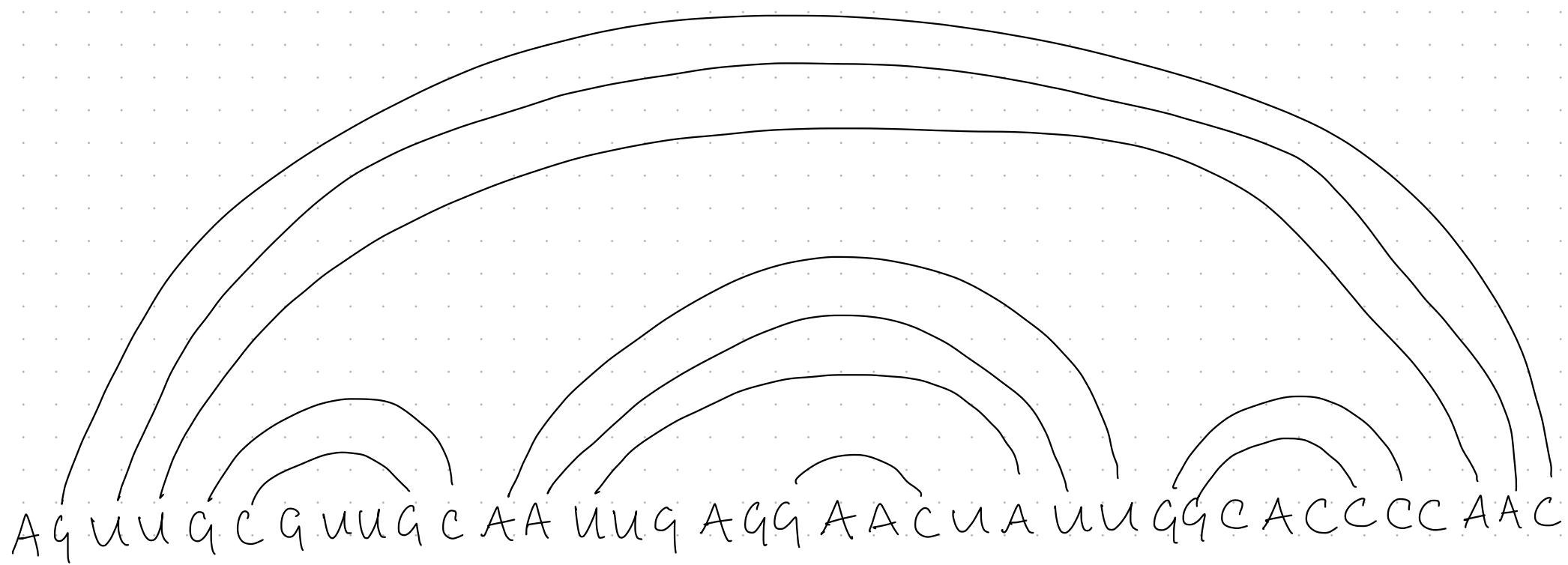
## Announcements

- \* HW 2 ongoing
- \* HW 1 grades released after lecture.
- \* Prelim 1      This Thursday      7:30-9 p.
  - \* Statler 185      } See Ed for room
  - \* Statler 196      } assignment.
- ↳ Coverage through last Friday's lecture.
- \* Review
  - \* Practice Exam on Canvas.
  - \* Review Session Tues      7-9 p
    - ↳ Gates G01.

Given. RNA Sequence me  $\{A, U, C, G\}^*$

Compute. Maximum Non-crossing matching E  
of preferred base pairs within m.

$$BP = \{(A, U), (C, G)\}$$



What are the cases / subproblems?

Cases In optimal folding of M.

\*

\*

M = uug c A A uug A G G A A C U A u u g g C A C C C C A A C

$M_1$

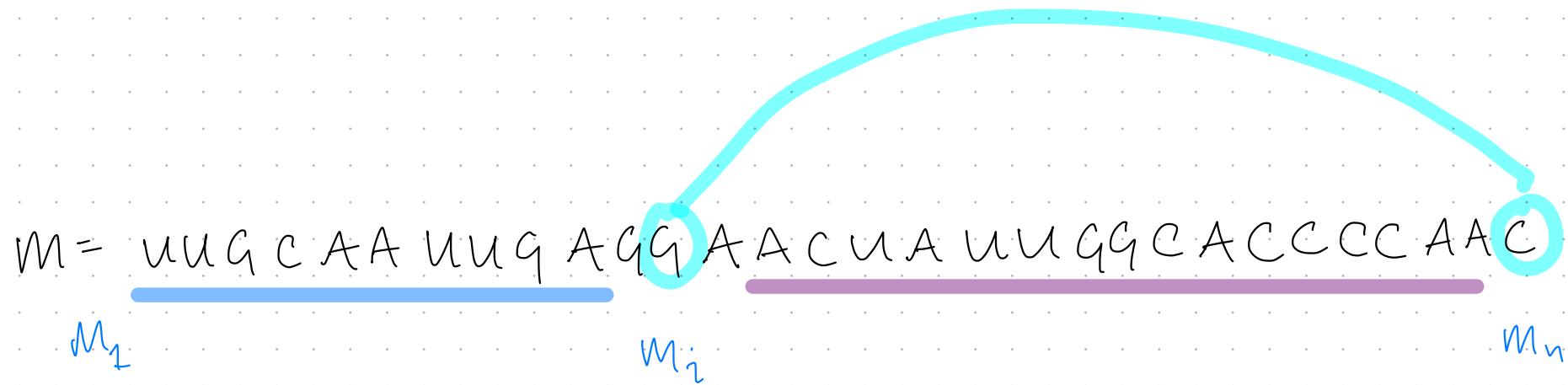
$M_n$

What happens to  $M_n$ ?

Cases In optimal folding of M.

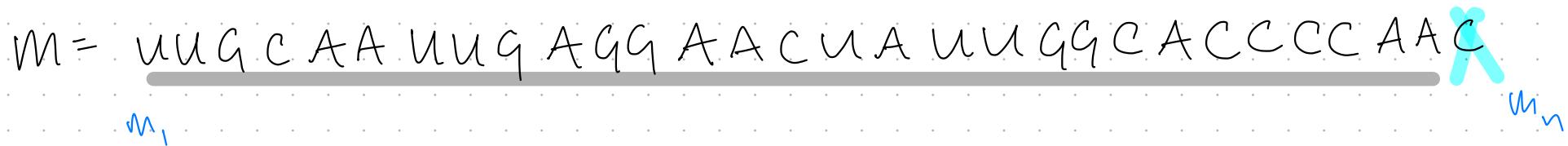
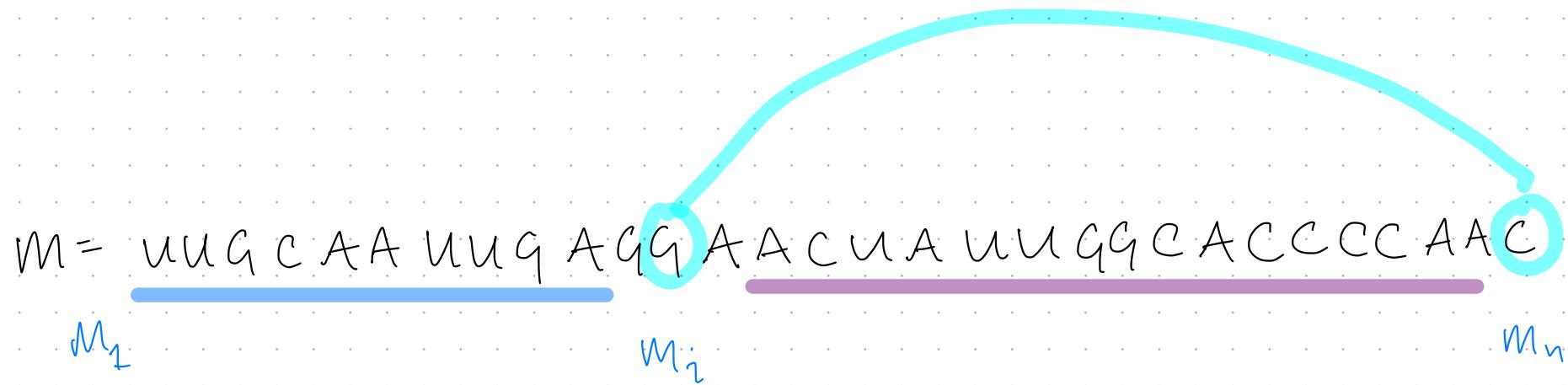
\*  $m_n$  base pairs w/ some  $m_i$   $i < n$ , OR

\*



## Cases In optimal folding of M.

- \*  $m_n$  base pairs w/ some  $m_i$   $i < n$ , OR
- \*  $m_n$  does NOT base pair



Cases

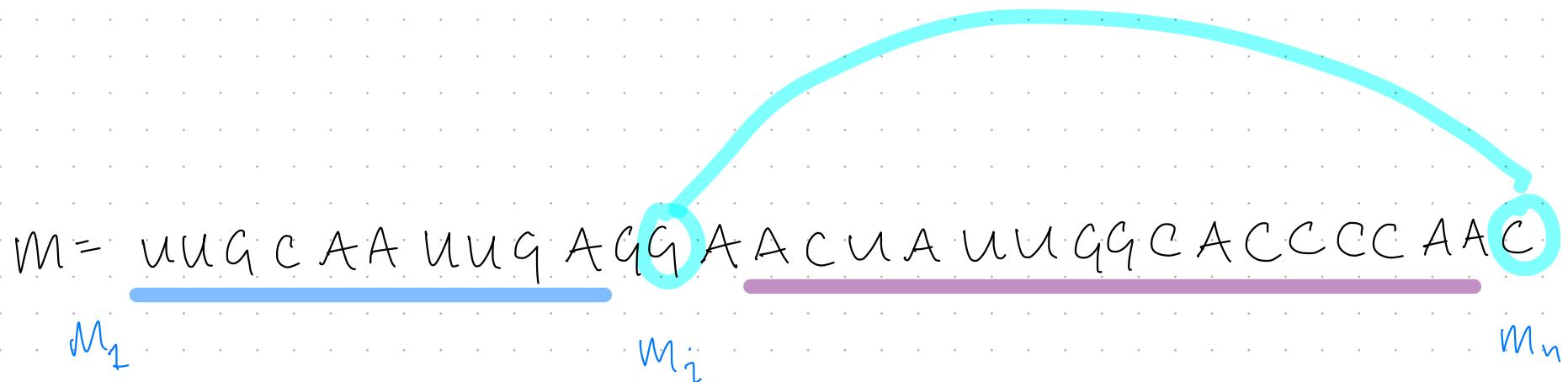
$M_n$  does Not base pair

$M = \text{UUG CAA UUG AGG A C U A U U G G C A C C C C A A C}$



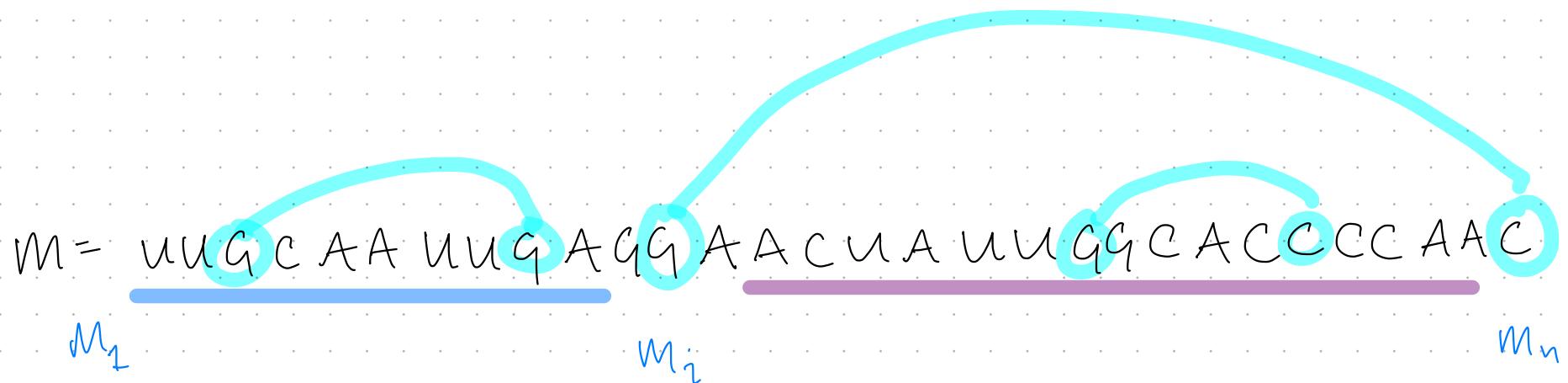
Best up til  $M_n = \text{Best up til } M_{n-1}$

Cases  $m_n$  base pairs w/ some  $m_i$   $i < n$



What base pairs are still legal?

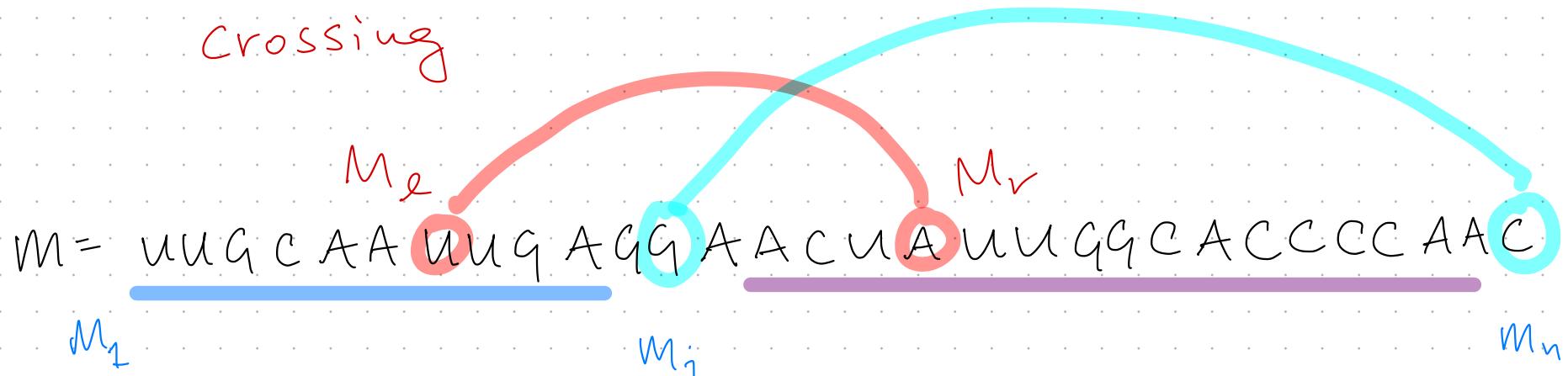
Cases  $m_n$  base pairs w/ some  $M_i$   $i < n$



What base pairs are still legal?

\* ALLOWED: from  $M_1 \rightarrow M_{i-1}$  &  $M_{i+1} \rightarrow M_{n-1}$

Cases  $m_n$  base pairs w/ some  $M_i$   $i < n$



What base pairs are still legal?

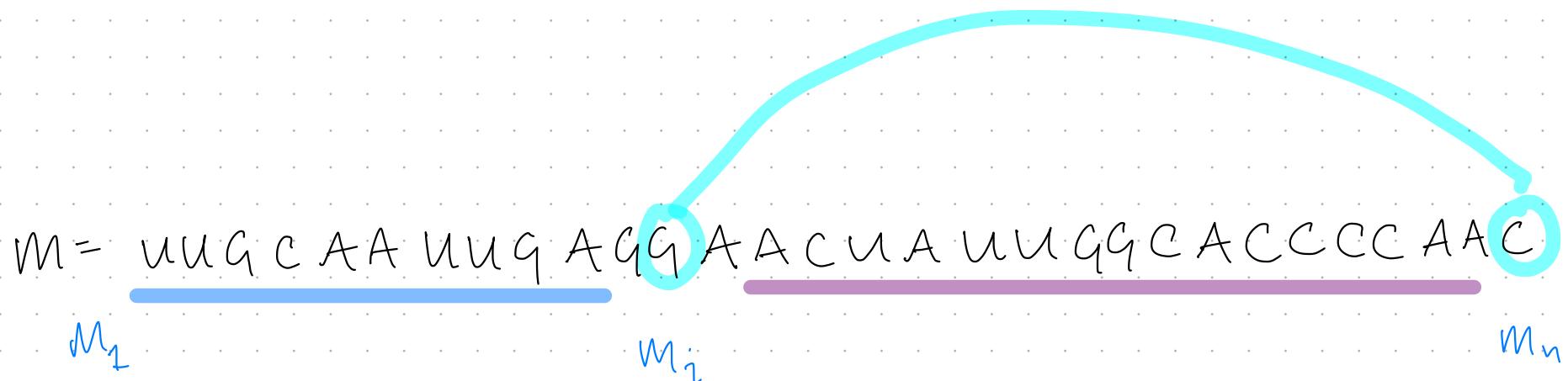
\* ALLOWED: from  $M_1 \rightarrow M_{i-1}$  &  $M_{i+1} \rightarrow M_{n-1}$

\* NOT ALLOWED: Crossings!

from  $M_\ell \rightarrow M_r$

for  $\ell < i < r$ .

Cases  $m_n$  base pairs w/ some  $m_i$   $i < n$



Best for  $M$

=

Best of left + Best of right

Subproblems = Substrings

$M = \text{UUG CAA UUG AGG AAC UAU UGG CAC CCC AAC}$

$M_l$

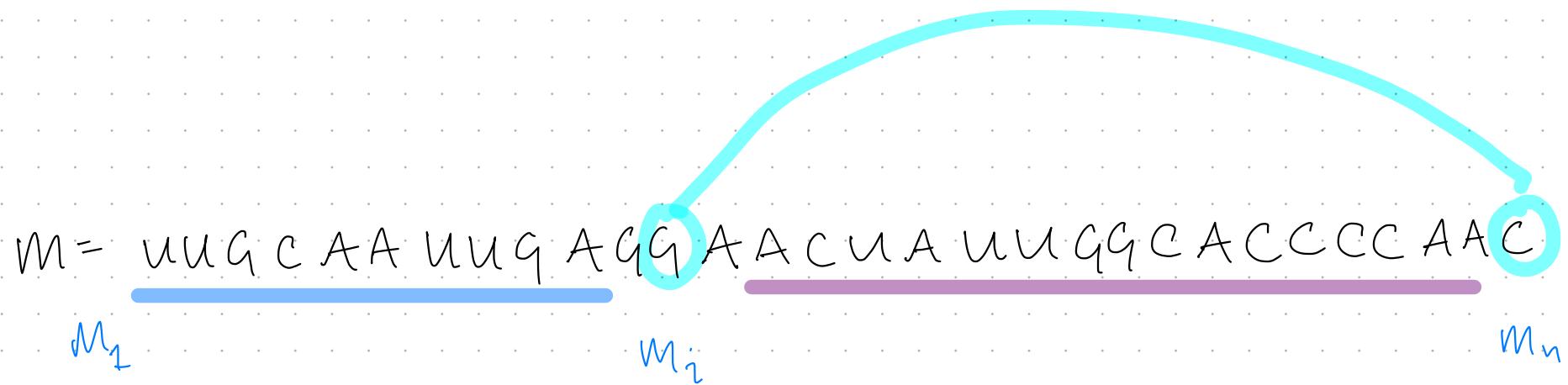
$M_r$

Best folding of substring from  
left endpoint  $M_l$  to right endpoint  $M_r$

= Fold( $l, r$ )

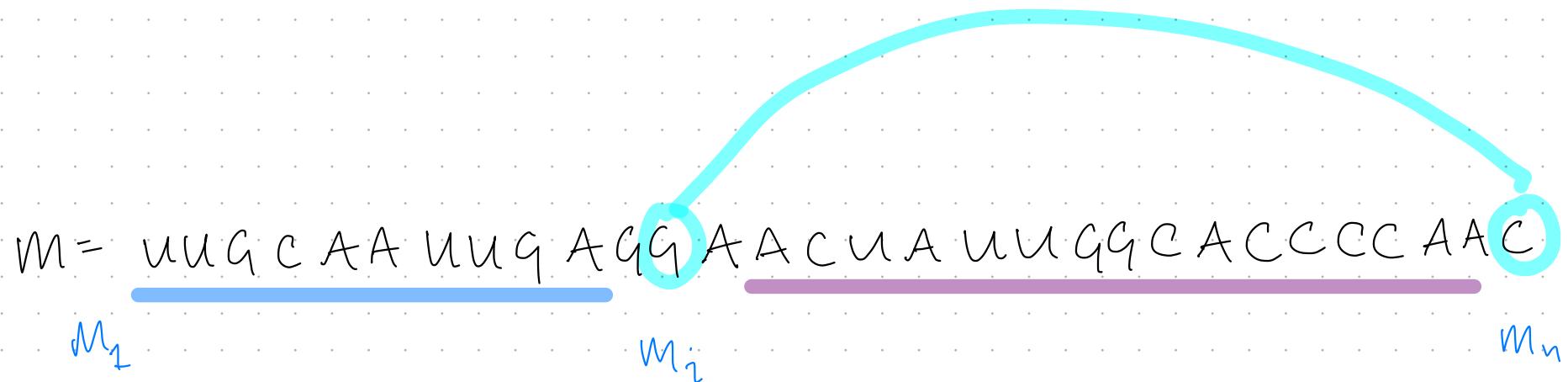
$\Rightarrow$  RNA Folding solved by  $\text{Fold}(1, n)$

Cases  $m_n$  base pairs w/ some  $m_i$   $i < n$



$$\text{Fold}(1, n) = 1 + \text{Fold}(1, i-1) + \text{Fold}(i+1, n-1)$$

Cases  $m_n$  base pairs w/ some  $m_i$   $i < n$



$$\text{Fold}(1, n) = \underline{1} + \underline{\text{Fold}(1, i-1)} + \underline{\text{Fold}(i+1, n-1)}$$

Note: subproblem does NOT start at  $l=1$

Cases

$m_n$  does Not base pair

$M = \text{UUGCAAUUGAGGACUAAUUGGACCCCCAAC}$



$m_1$   $m_n$

$$\text{Fold}(1, n) = \text{Fold}(1, n-1)$$

$M = \underline{\text{uug caa uug agg aac ua uuggc acccc aac}}$

$M_1$

$M_i$

$M_n$

$M = \underline{\text{uug caa uug agg aac ua uuggc acccc aac}}$

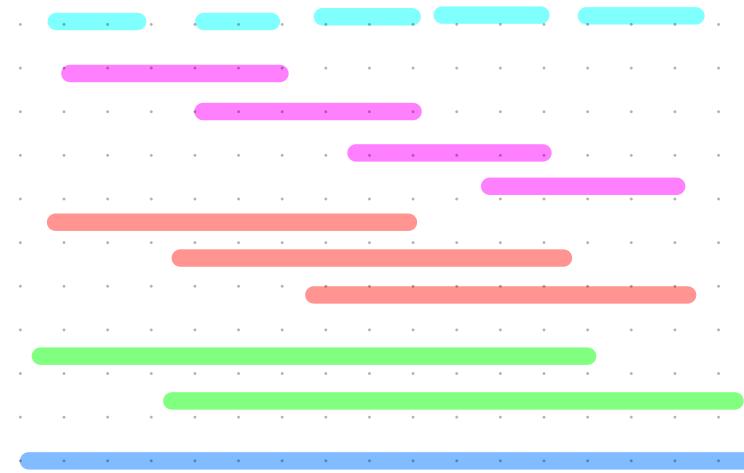
$m_1$

$m_n$

## RNA Recurrence

$$\text{Fold}(1, n) = \max \left\{ \begin{array}{l} 1 + \max_{i < n} \left\{ \text{Fold}(1, i-1) + \text{Fold}(i+1, n) \right. \\ \text{s.t. } (M_i, M_n) \in \text{BP} \\ \text{Fold}(1, n-1) \end{array} \right\}$$

$M = A C A G U$



## Subproblems

\* For each width

1, 2, 3, 4, 5

\* For each left endpoint

A, C, A, G, U

A

AC

ACA

ACAG

ACAGU

C

CA

CAG

A

AG

G

A

AC

ACA

ACAG

ACAGIX

C

CA

CAG

A

AG

G

A

AC

ACA

ACAG

ACAGU

C

CA

CAG

A

AG

G

A

AC

ACA

ACAG

ACAGU

C

CA

CAG

A

AG

G

A

A X

ACA X

ACAX X

ACAGIX

C

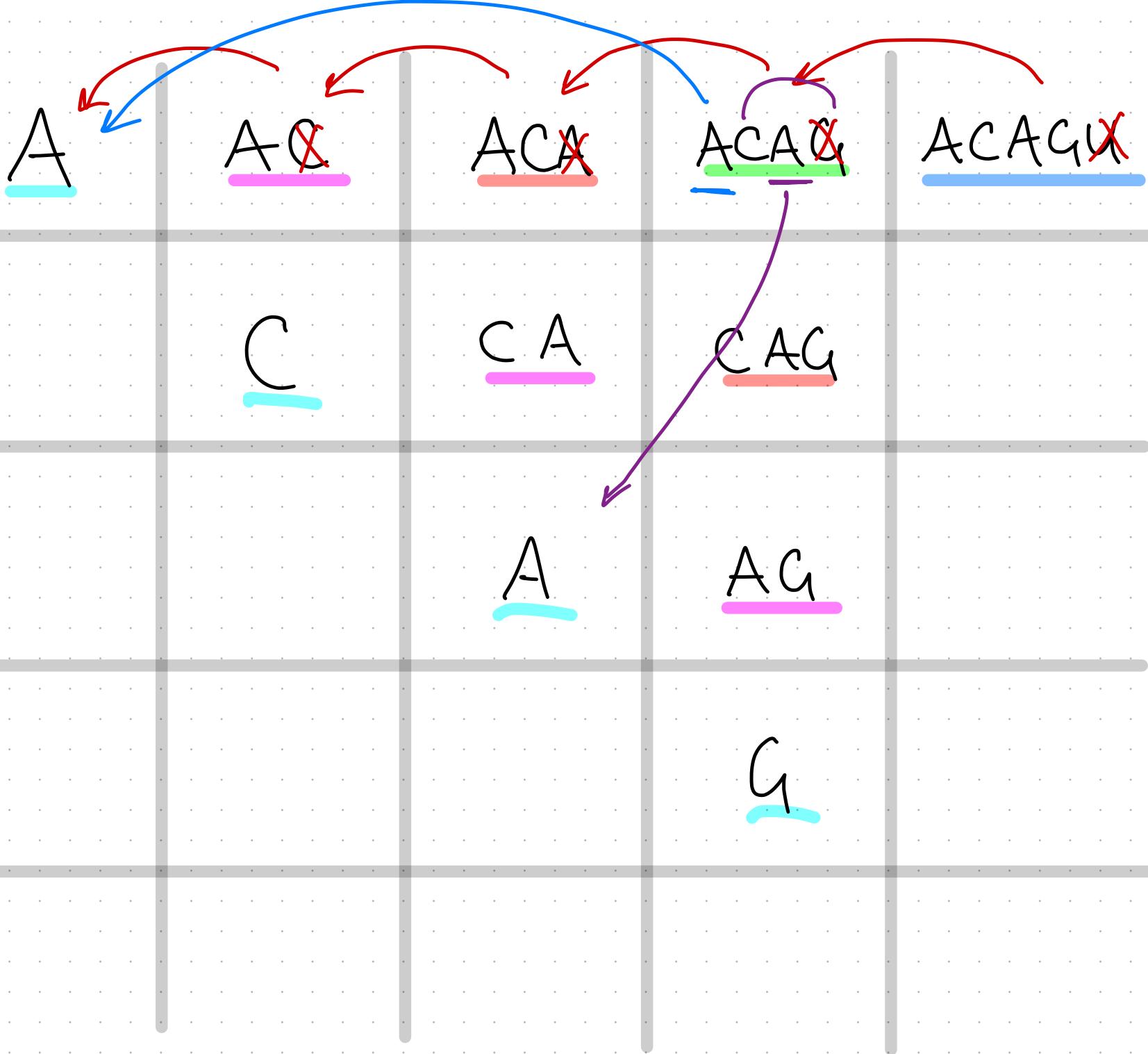
CA

CAG

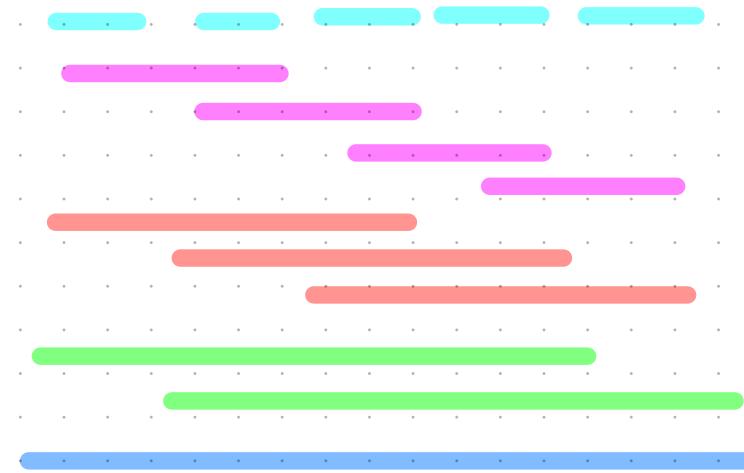
A

AG

G



$M = A C A G U$



## Subproblems,

\* For each width

1, 2, 3, 4, 5

\* For each left endpoint

A, C, A, G, U

## RNA-Fold(m)

For  $l = 1 \rightarrow n$  :  $\text{Fold}(l, l) = 0$ .

For width =  $1 \rightarrow n-1$

For  $l = 1 \rightarrow n - \text{width}$

$$r = l + \text{width}$$

no Pair  $\leftarrow \text{Fold}(l, r-1)$

best Pair  $\leftarrow$

Find best legal base pair  
between  $M_l$  and  $M_r$

for some  $l < q < r$

$$\text{Fold}(l, r) = \max \left\{ \text{no Pair}, \text{best Pair} \right\}$$

## RNA-Fold(m)

For  $l = 1 \rightarrow n$  :  $\text{Fold}(l, l) = 0$ .

For width =  $1 \rightarrow n-1$

For  $l = 1 \rightarrow n - \text{width}$

$$r = l + \text{width}$$

no Pair  $\leftarrow \text{Fold}(l, r-1)$

best Pair  $\leftarrow 0$

For  $q = l \rightarrow r-1$

if  $(M_q, M_r) \in \text{BP}$

qPair  $\leftarrow 1 + \text{Fold}(l, q-1) + \text{Fold}(q+1, r-1)$

bestPair  $\leftarrow \max \{ q\text{Pair}, \text{bestPair} \}$

$\text{Fold}(l, r) = \max \{ \text{no Pair}, \text{best Pair} \}$

A

Fold (1,1)

C

Fold (2,2)

A

Fold (3,3)

G

Fold (4,4)

A



AC ~~X~~

Fold(1,1)

Fold(1,2)

C

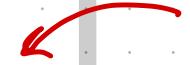


CA ~~X~~

Fold(2,2)

Fold(2,3)

A



AA ~~X~~

Fold(3,3)

Fold(3,4)

G

Fold(4,4)

A

A ~~C~~

AC ~~A~~

Fold(1,1)

Fold(1,2)

Fold(1,3)

C

CA ~~A~~

CA ~~G~~

Fold(2,2)

Fold(2,3)

Fold(2,4)

A

AC ~~A~~

Fold(3,3)

Fold(3,4)

G

Fold(4,4)

A

AC

ACA

ACAG

Fold(1,1)

Fold(1,2)

Fold(1,3)

Fold(1,4)

C

CA

CAG

Fold(2,2)

Fold(2,3)

Fold(2,4)

A

AG

Fold(3,3)

Fold(3,4)

G

Fold(4,4)

A

Fold(1,1)

AC~~X~~

Fold(1,2)

ACA~~X~~

Fold(1,3)

ACAG~~X~~

Fold(1,4)

ACAGU~~X~~

Fold(1,5)

C

Fold(2,2)

CA~~X~~

Fold(2,3)

CAG~~X~~

Fold(2,4)

A

Fold(3,3)

AG~~X~~

Fold(3,4)

G

Fold(4,4)

## RNA-Fold(m)

For  $l = 1 \rightarrow n$  :  $\text{Fold}(l, l) = 0$ .

For width =  $1 \rightarrow n-1$

For  $l = 1 \rightarrow n - \text{width}$

$$r = l + \text{width}$$

no Pair  $\leftarrow \text{Fold}(l, r-1)$

best Pair  $\leftarrow 0$

For  $q = l \rightarrow r-1$

if  $(M_q, M_r) \in \text{BP}$

qPair  $\leftarrow 1 + \text{Fold}(l, q-1) + \text{Fold}(q+1, r-1)$

bestPair  $\leftarrow \max \{ q\text{Pair}, \text{bestPair} \}$

$\text{Fold}(l, r) = \max \{ \text{no Pair}, \text{best Pair} \}$

Running Time

Entries for each  
 $(l, r) \in [n] \times [n]$

$\Rightarrow O(n^2)$

To fill in each entry

↳ Probe q s.t.

$$l < q < r$$

$O(n)$

Total

$O(n^3)$

## Proof of Correctness.

For all  $n \in \mathbb{N}$ ,  $\text{Fold}(1, n)$  returns max non-crossing match.

## Proof of Correctness.

For all  $n \in \mathbb{N}$ ,  $\text{Fold}(1, n)$  returns max non-crossing match.



For all widths  $w \in \mathbb{N}$ , for all left endpoints  $l \leq |M| - w$

$\text{Fold}(l, l+w) = \max$  non-crossing match  
w/ left endpoint  $M_l$   
and right endpoint  $M_{l+w}$ .

## Proof of Correctness.

For all  $n \in \mathbb{N}$ ,  $\text{Fold}(1, n)$  returns max non-crossing match.



For all widths  $w \in \mathbb{N}$ , for all left endpoints  $l \leq |M| - w$

$\text{Fold}(l, l+w) = \max \text{ non-crossing match}$   
w/ left endpoint  $M_l$   
and right endpoint  $M_{l+w}$ .

---

Base case.  $w=0$ .

\* No base pairing possible w/ single character  
 $\Rightarrow \max \text{ match} = 0$   
for all  $l$ .

\*  $\text{Fold}(l, l) = 0$



Inductive Hypothesis,  $\forall w_0 < w$ , for all legal  $\ell$

$\text{Fold}(\ell, \ell + w_0)$  = best RNA fold from  $M_\ell$  to  $M_{\ell+w_0}$

Consider

$M_\ell M_{\ell+1} M_{\ell+2} \dots M_{\ell+w-1} M_{\ell+w}$ .

Two cases.

\*  $M_{\ell+w}$  not paired :

Best fold from  $M_\ell \rightarrow M_{\ell+w}$

= Best fold from  $M_\ell \rightarrow M_{\ell+w-1}$ ) By IH

=  $\text{Fold}(\ell, \ell + w - 1)$



width  $w-1 < w$

Inductive Hypothesis,  $\forall w_0 < w$ , for all legal  $\ell$

$\text{Fold}(\ell, \ell + w_0)$  = best RNA fold from  $M_\ell$  to  $M_{\ell+w_0}$

Consider

$M_\ell M_{\ell+1} M_{\ell+2} \cdots M_{\ell+w-1} M_{\ell+w}$ .

Two cases.  $M_{\ell+w}$  paired :

Best fold from  $M_\ell \rightarrow M_{\ell+w}$

$$= \text{Best of } \text{BP}(M_q, M_r) +$$

Best fold from  $M_\ell \rightarrow M_{q-1}$  + Best fold from

$M_{q+1} \rightarrow M_{\ell+w-1}$

Inductive Hypothesis,  $\forall w_0 < w$ , for all legal  $\ell$

$\text{Fold}(\ell, \ell + w_0)$  = best RNA fold from  $M_\ell$  to  $M_{\ell+w_0}$

Consider

$M_\ell M_{\ell+1} M_{\ell+2} \cdots M_{\ell+w-1} M_{\ell+w}$ .

Two cases.  $M_{\ell+w}$  paired :

Best fold from  $M_\ell \rightarrow M_{\ell+w}$

$$= \text{Best of } \text{BP}(M_q, M_r) +$$

Both of width  
 $w' < w$ .

Best fold from  $M_\ell \rightarrow M_{q-1}$  + Best fold from

$M_{q+1} \rightarrow M_{\ell+w-1}$

Inductive Hypothesis,  $\forall w_0 < w$ , for all legal  $\ell$

$\text{Fold}(\ell, \ell + w_0)$  = best RNA fold from  $M_\ell$  to  $M_{\ell+w_0}$

Consider

$M_\ell M_{\ell+1} M_{\ell+2} \cdots M_{\ell+w-1} M_{\ell+w}$ .

Two cases.  $M_{\ell+w}$  paired :

Best fold from  $M_\ell \rightarrow M_{\ell+w}$

$$= \text{Best of } \text{BP}(M_q, M_r) +$$

Both of width  
 $w' < w$ .

Best fold from  $M_\ell \rightarrow M_{q-1}$  + Best fold from

$M_{q+1} \rightarrow M_{\ell+w-1}$

$$\approx 1 + \text{Fold}(\ell, q-1) + \text{Fold}(q+1, \ell+w-1) \quad \checkmark \text{ By IH}$$

