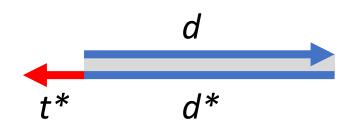# DNA sequence design

slides © 2021, David Doty
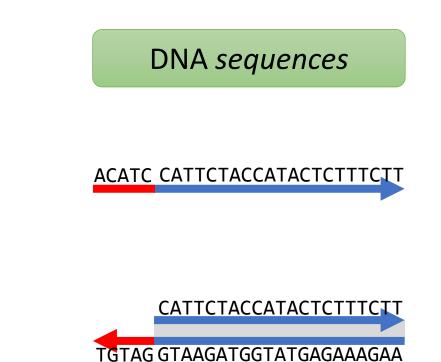
ECS 232: Theory of Molecular Computation, UC Davis

# Two layers of abstraction in DNA nanotech

DNA *strands* with abstract "binding domains"
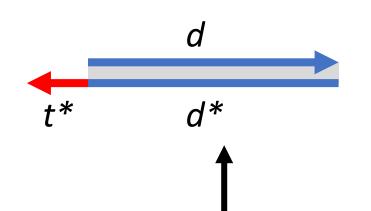
DNA *sequences*

# Two layers of abstraction in DNA nanotech

DNA *strands* with abstract "binding domains"

DNA *sequences*



This describes ideally how we **want** strands to bind.

2

# Two layers of abstraction in DNA nanotech

# DNA sequence design



bad choice of DNA sequence

$t$  $d$

GGCCG GCCGGTTTTTCCGGCCGGCCAAT

# DNA sequence design



bad choice of DNA sequence

$t$      $d$

GGCCG GCCGGTTTTTCCGGCCGGCCAAT

most likely structure

CCGGCCGGCCAAT
GGCCGGCCGG

# DNA sequence design



bad choice of
DNA sequence

*t*     *d*

GGCCG GCCGGTTTTTCCGGCCGGCCAAT

most likely structure

CCGGCCGGCCAAT
GGCCGGCCGG

*Why is this bad*?
If we want the strand to bind to other strands,
it first has to <u>break up</u> its own structure.
i.e., *effective* binding rate/strength is lowered

3

# Common DNA sequence design goals: <u>What to avoid</u>

- Excessive secondary structure of strands

# Common DNA sequence design goals: <u>What to avoid</u>

- Excessive secondary structure of strands

  CCGGCCGGCCAAT

  GGCCGGCCGG

- Significant interaction between non-complementary domains

  domain $d$

  AAAAAAAAAAAAAAAAAAAAAA

  TTTTTTGTTTTTTTTATTTT          TTTTTTTTTTTTTTTTTTTTT

  domain $f*$                         domain $d*$

# Common DNA sequence design goals: <u>What to avoid</u>

- Excessive secondary structure of strands

CCGGCCGGCCAAT
GGCCGGCCGG

- Significant interaction between non-complementary domains

domain *d*

AAAAAAAAAAAAAAAAAAAAA

TTTTTTGTTTTTTTTATTTT          TTTTTTTTTTTTTTTTTTTTT

domain *f\**                    domain *d\**

- Certain string-based rules, e.g.
  - some patterns such as GGGG (forms "*G-tetraplex*": https://www.idtdna.com/pages/education/decoded/article/g-repeats-structural-challenges-for-oligo-design)
  - > 70 % or < 30% G/C content (G/C binds more strongly)
  - domains ending in A/T (they "breathe" more)

- And often other constraints

# DNA energy models

How do we predict what structures DNA strands are likely to form?

# DNA <u>duplex</u> energy model (simple versions)

- How strongly does a DNA strand bind to its <u>perfect complement</u>?

# DNA <u>duplex</u> energy model (simple versions)

- How strongly does a DNA strand bind to its <u>perfect complement</u>?
- 1<sup>st</sup> approximation: <span style="color:red">proportional to length</span>:
  - $\Delta G$( 5'-AAGGTTAC-3' ,

    3'-TTCCAATG-5' ) = 1+1+1+1+1+1+1+1 = 8

# DNA <u>duplex</u> energy model (simple versions)

- How strongly does a DNA strand bind to its <u>perfect complement</u>?

- <u>1$^{st}$ approximation</u>: <span style="color:red">proportional to length</span>:
  - $\Delta G$( 5'-AAGGTTAC-3' ,

      3'-TTCCAATG-5' ) = 1+1+1+1+1+1+1+1 = 8

- <u>2$^{nd}$ approximation</u>: <span style="color:red">depends on base pair</span>:
  - G/C about twice as strong as A/T
  - $\Delta G$( 5'-AAGGTTAC-3' ,

      3'-TTCCAATG-5' ) = 1+1+2+2+1+1+1+2 = 11

# DNA <u>duplex</u> energy model (simple versions)

- How strongly does a DNA strand bind to its <u>perfect complement</u>?

- 1<sup>st</sup> <u>approximation</u>: <span style="color:red">proportional to length</span>:
  - $\Delta G$( 5'-AAGGTTAC-3' ,
    3'-TTCCAATG-5' ) = 1+1+1+1+1+1+1+1 = 8

- 2<sup>nd</sup> <u>approximation</u>: <span style="color:red">depends on base pair</span>:
  - G/C about twice as strong as A/T
  - $\Delta G$( 5'-AAGGTTAC-3' ,
    3'-TTCCAATG-5' ) = 1+1+2+2+1+1+1+2 = 11

- 3<sup>rd</sup> <u>approximation</u>: <span style="color:red">nearest neighbor model</span> (used in practice):
  - depends on base pair, *and* on the neighboring base pairs

# Why do the neighbors matter?

Responsible for specific base-pairing

Holds adjacent bases together on a single strand (strong covalent bond)

Much of DNA stability is not from base pair (formed by hydrogen bonds) but from "stacking" interactions between adjacent bases.



hydrogen bonding

stacking

backbone

coaxial stacking

cross stacking

source: https://dna-robotics.eu/2019/11/29/simulating-dna/

7

# Nearest neighbor energy model

$$\downarrow \quad \downarrow \quad \downarrow$$

$$5' \; C\text{-}G\text{-}T\text{-}T\text{-}G\text{-}A \; 3'$$
$$* \; * \; * \; * \; * \; *$$
$$3' \; G\text{-}C\text{-}A\text{-}A\text{-}C\text{-}T \; 5'$$
$$\uparrow \qquad \uparrow$$

$\Delta G^\circ_{37}(\text{pred.}) = \Delta G^\circ(\text{CG/GC}) + \Delta G^\circ(\text{GT/CA}) + \Delta G^\circ(\text{TT/AA})$

$+ \Delta G^\circ(\text{TG/AC}) + \Delta G^\circ(\text{GA/CT}) + \Delta G^\circ(\text{init.})$

$= -2.17 - 1.44 - 1.00 - 1.45 - 1.30 + \boxed{0.98 + 1.03}$

$\Delta G^\circ_{37}(\text{pred.}) = -5.35$ kcal/mol

$\Delta G^\circ_{37}(\text{obs.}) = -5.20$ kcal/mol

$\Delta G_{\text{init}}$ = penalty for bringing together two strands (TODO: maybe not... not explained in paper) (*different terms if end is C/G or A/T*)

[*A unified view of polymer, dumbbell, and oligonucleotide DNA nearest-neighbor thermodynamics,* John SantaLucia Jr., PNAS 1998]

| Table 1. Compari | |
| --- | --- |
| Sequence | Unified (ref. 22) |
| AA/TT | −1.00 |
| AT/TA | −0.88 |
| TA/AT | −0.58 |
| CA/GT | −1.45 |
| GT/CA | −1.44 |
| CT/GA | −1.28 |
| GA/CT | −1.30 |
| CG/GC | −2.17 |
| GC/CG | −2.24 |
| GG/CC | −1.84 |
| Average | −1.42 |

# Energy of <u>non-duplex</u> secondary structures

What about DNA strands that are not perfectly complementary, but *some* bases match?

# Energy of **non-duplex** secondary structures

What about DNA strands that are not perfectly complementary, but *some* bases match?

**Definition**: A secondary structure of a set of DNA strands is a set of base pairs among complementary bases.
Formally, it is a *matching* on the graph $G=(V,E)$, where
$V = \{$ bases in each strand $\}$
$E = \{ \{u,v\} \mid \{u,v\} = \{A,T\}$ or $\{u,v\} = \{G,C\} \}$

# Energy of <u>non-duplex</u> secondary structures

What about DNA strands that are not perfectly complementary, but *some* bases match?

unpseudoknotted:

**Definition**: A <u>secondary structure</u> of a set of DNA strands is a set of base pairs among complementary bases. Formally, it is a *matching* on the graph $G=(V,E)$, where
$V$ = { bases in each strand }
$E$ = { {$u,v$} | {$u,v$} = {A,T} or {$u,v$} = {G,C} }

sometimes drawn with strands straight and base pairs as curved arcs:

**Definition**: A secondary structure is <u>unpseudoknotted</u> (with respect to a particular circular ordering of the strands) if, drawing strands in 5'-3' order in a *circle* and connecting the base pairs by *straight lines*, **no lines cross**.

GCAGGAGUCUAGCGAUGCUAGUCAGCUAGCUCAUAAUGAAUAGGCUACGACUAGCGCUGGAGACCUU

# Energy of **non-duplex** secondary structures

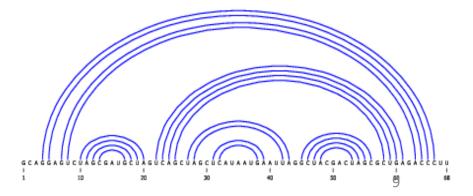What about DNA strands that are not perfectly complementary, but *some* bases match?

**Definition**: A secondary structure of a set of DNA strands is a set of base pairs among complementary bases. Formally, it is a *matching* on the graph $G=(V,E)$, where
$V = \{$ bases in each strand $\}$
$E = \{ \{u,v\} \mid \{u,v\} = \{A,T\} \text{ or } \{u,v\} = \{G,C\} \}$

unpseudoknotted:



sometimes drawn with strands straight and base pairs as curved arcs:
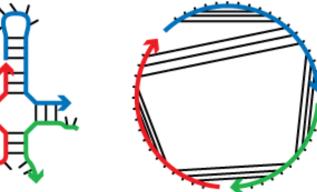
pseudoknots:

**Definition**: A secondary structure is unpseudoknotted (with respect to a particular circular ordering of the strands) if, drawing strands in 5'-3' order in a *circle* and connecting the base pairs by *straight lines*, **no lines cross**.





9

# Equivalent definitions of unpseudoknotted

**Definition 1**: Drawing strands in 5'-3' order in a *circle* and connecting the base pairs by *straight lines*, no lines cross.

# Equivalent definitions of unpseudoknotted

**Definition 1**: Drawing strands in 5'-3' order in a *circle* and connecting the base pairs by *straight lines*, no lines cross.



**Definition 2**: Base pair indices obey the **nesting property**: there are *no* base pairs $(a,b) \in \mathbb{N}^2$ and $(x,y) \in \mathbb{N}^2$ such that $a < x < b < y$ (e.g., it can be $a < b < x < y$ or $a < x < y < b$ )

# Equivalent definitions of unpseudoknotted

**Definition 1**: Drawing strands in 5'-3' order in a *circle* and connecting the base pairs by *straight lines*, no lines cross.



**Definition 2**: Base pair indices obey the **nesting property**: there are *no* base pairs $(a,b) \in \mathbb{N}^2$ and $(x,y) \in \mathbb{N}^2$ such that $a < x < b < y$ (e.g., it can be $a < b < x < y$ or $a < x < y < b$ )

**Definition 3**: Balanced parentheses describe base pairs in **dot-parens** (a.k.a., **dot-bracket**) notation.



5'         3'

$(((.....)))....(((.....)))$.

# Equivalent definitions of unpseudoknotted

**Definition 1**: Drawing strands in 5'-3' order in a *circle* and connecting the base pairs by *straight lines*, no lines cross.



**Definition 2**: Base pair indices obey the **nesting property**: there are *no* base pairs $(a,b) \in \mathbb{N}^2$ and $(x,y) \in \mathbb{N}^2$ such that $a < x < b < y$ (e.g., it can be $a < b < x < y$ or $a < x < y < b$ )

**Definition 3**: Balanced parentheses describe base pairs in **dot-parens** (a.k.a., **dot-bracket**) notation.



pseudoknotted: need multiple parenthesis types to describe

$(((.....)))....(((.....)))$.

$(((......[[)))...]]]..$
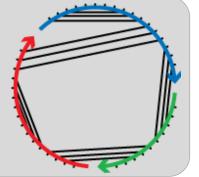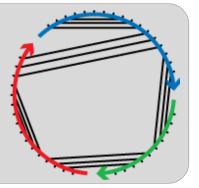
# Equivalent definitions of unpseudoknotted

**Definition 1**: Drawing strands in 5'-3' order in a *circle* and connecting the base pairs by *straight lines*, no lines cross.
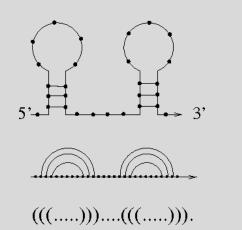
**Definition 2**: Base pair indices obey the **nesting property**: there are *no* base pairs $(a,b) \in \mathbb{N}^2$ and $(x,y) \in \mathbb{N}^2$ such that $a < x < b < y$   (e.g., it can be $a < b < x < y$ or $a < x < y < b$ )

**Definition 3**: Balanced parentheses describe base pairs in **dot-parens** (a.k.a., **dot-bracket**) notation.

pseudoknotted: need multiple parenthesis types to describe

$(((.....)))....(((.....))).$

$(((......[[)))...]]]..$

**Definition 4**: The graph $G=(V,E)$ is **outerplanar**, where
$V$ = { bases in each strand }
$E$ = { {$u,v$} |    {$u,v$} are a paired base pair,
          or {$u,v$} are adjacent }

outerplanar = can be drawn with no edges crossing (planar), **and** all vertices incident to the outer face

outerplanar          not outerplanar

# Back to first approximation of energy model

- (For now, consider only one strand.)
- Given a DNA sequence $S$, what is the maximum number of base pairs that can be formed in any unpseudoknotted secondary structure?

# Back to first approximation of energy model

- (For now, consider only one strand.)
- Given a DNA sequence $S$, what is the maximum number of base pairs that can be formed in any unpseudoknotted secondary structure?
  - Without unpseudoknotted constraint, this is trivial:   min(#C,#G) + min(#A,#T)

# Back to first approximation of energy model

- (For now, consider only one strand.)
- Given a DNA sequence *S*, what is the maximum number of base pairs that can be formed in any unpseudoknotted secondary structure?
  - Without unpseudoknotted constraint, this is trivial:   min(#C,#G) + min(#A,#T)
- Can be taken as a rough approximation of the <span style="color:red">minimum free energy</span> structure of *S*, i.e., the <span style="color:red">most probable</span> structure "at thermodynamic equilibrium" (*what you'd see if you heat it up and slowly cool it*).

# Computing maximally bound unpseudoknotted secondary structure in polynomial time



1   2            $i$                                    $j$                        $n$

**Recursive solution**:
- Strand length is $n$.
- For $1 \leq i \leq j \leq n$, let OPT($i,j$) = max base pairs possible using **only** bases $i$ through $j$.

# Computing maximally bound unpseudoknotted secondary structure in polynomial time

pair $j$ with another base or not?



1  2          $i$                                    $j–1$  $j$                              $n$

**Recursive solution**:
- Strand length is $n$.
- For $1 \leq i \leq j \leq n$, let OPT$(i,j)$ = max base pairs possible using **only** bases $i$ through $j$.
- <u>Question</u>: do we pair base $j$ with some other base between $i$ and $j–1$?

# Computing maximally bound unpseudoknotted secondary structure in polynomial time

pair *j* with another base or not?



1  2        *i*                                                    *j*–1  *j*                    *n*

**Recursive solution**:
- Strand length is *n*.
- For $1 \le i \le j \le n$, let OPT(*i*,*j*) = max base pairs possible using **only** bases *i* through *j*.
- Question: do we pair base *j* with some other base between *i* and *j*–1?
- If *not*, recursively, the optimal value is:
    - OPT(*i*,*j*) = OPT(*i*,*j*–1)

# Computing maximally bound unpseudoknotted secondary structure in polynomial time



pair $j$ with another base or not?

| 1 | 2 | | $i$ | | | | $k-1$ | $k$ | $k+1$ | | | | $j-1$ | $j$ | | | | | | $n$ |

**Recursive solution**:
- Strand length is $n$.
- For $1 \leq i \leq j \leq n$, let OPT($i,j$) = max base pairs possible using **only** bases $i$ through $j$.
- Question: do we pair base $j$ with some other base between $i$ and $j-1$?
- If *not*, recursively, the optimal value is:
  - OPT($i,j$) = OPT($i,j-1$)
- If we pair $j$ with $k$, **nesting property** implies no base pair can form between any base in [$i,\dots k-1$] and any base in [$k+1,j-1$]

# Computing maximally bound unpseudoknotted secondary structure in polynomial time

pair $j$ with another base or not?



1   2                 $i$                        $k{-}1$   $k$   $k{+}1$                        $j{-}1$   $j$                                    $n$

**Recursive solution**:
- Strand length is $n$.
- For $1 \leq i \leq j \leq n$, let OPT($i,j$) = max base pairs possible using **only** bases $i$ through $j$.
- <u>Question</u>: do we pair base $j$ with some other base between $i$ and $j{-}1$?
- If *not*, recursively, the optimal value is:
  - OPT($i,j$) = OPT($i,j{-}1$)
- If we pair $j$ with $k$, **nesting property** implies no base pair can form between any base in $[i,\dots k{-}1]$ and any base in $[k{+}1,j{-}1]$
- Recursively, optimal value depends on:
  - OPT($i,k{-}1$) and OPT($k{+}1,j{-}1$)

# Computing maximally bound unpseudoknotted secondary structure in polynomial time

pair $j$ with another base or not?



1  2     $i$     $k{-}1$  $k$  $k{+}1$     $j{-}1$  $j$     $n$

**Recursive solution**:
- Strand length is $n$.
- For $1 \le i \le j \le n$, let OPT$(i,j)$ = max base pairs possible using **only** bases $i$ through $j$.
- Question: do we pair base $j$ with some other base between $i$ and $j{-}1$?
- If *not*, recursively, the optimal value is:
  - OPT$(i,j)$ = OPT$(i,j{-}1)$
- If we pair $j$ with $k$, **nesting property** implies no base pair can form between any base in $[i,\ldots k{-}1]$ and any base in $[k{+}1,j{-}1]$
- Recursively, optimal value depends on:
  - OPT$(i,k{-}1)$ and OPT$(k{+}1,j{-}1)$

**Recursive algorithm** (implement w/ dynamic programming):
OPT$(i,j)$ = max of:

12

# Computing maximally bound unpseudoknotted secondary structure in polynomial time



pair $j$ with another base or not?

1   2        $i$           $k{-}1$  $k$  $k{+}1$              $j{-}1$  $j$           $n$
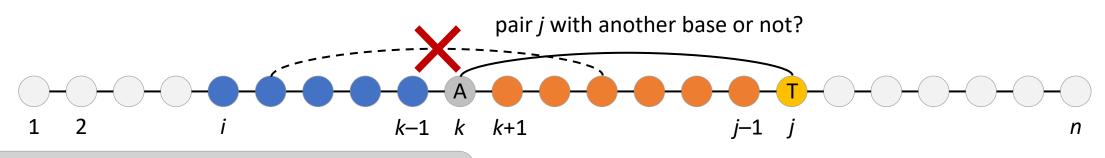
**Recursive solution**:
- Strand length is $n$.
- For $1 \le i \le j \le n$, let OPT$(i,j)$ = max base pairs possible using **only** bases $i$ through $j$.
- Question: do we pair base $j$ with some other base between $i$ and $j{-}1$?
- If *not*, recursively, the optimal value is:
  - OPT$(i,j)$ = OPT$(i,j{-}1)$
- If we pair $j$ with $k$, **nesting property** implies no base pair can form between any base in $[i,\dots k{-}1]$ and any base in $[k{+}1,j{-}1]$
- Recursively, optimal value depends on:
  - OPT$(i,k{-}1)$ and OPT$(k{+}1,j{-}1)$

**Recursive algorithm** (implement w/ dynamic programming)**:**
OPT$(i,j)$ = max of:
   OPT$(i,j{-}1)$,       *// don't form base pair with $j$*

12

# Computing maximally bound unpseudoknotted secondary structure in polynomial time



pair $j$ with another base or not?

1   2        $i$              $k-1$   $k$   $k+1$            $j-1$   $j$                n

**Recursive solution:**
- Strand length is $n$.
- For $1 \leq i \leq j \leq n$, let OPT($i,j$) = max base pairs possible using **only** bases $i$ through $j$.
- <u>Question</u>: do we pair base $j$ with some other base between $i$ and $j-1$?
- If *not*, recursively, the optimal value is:
  - OPT($i,j$) = OPT($i,j-1$)
- If we pair $j$ with $k$, **nesting property** implies no base pair can form between any base in $[i,... k-1]$ and any base in $[k+1,j-1]$
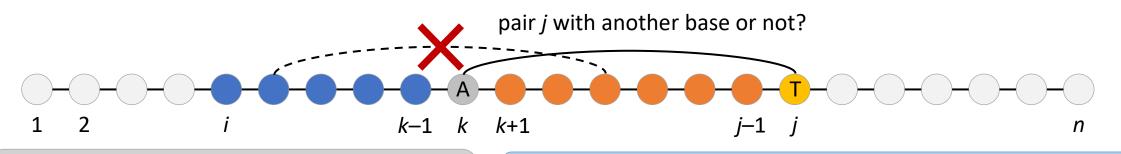- Recursively, optimal value depends on:
  - OPT($i,k-1$) and OPT($k+1,j-1$)

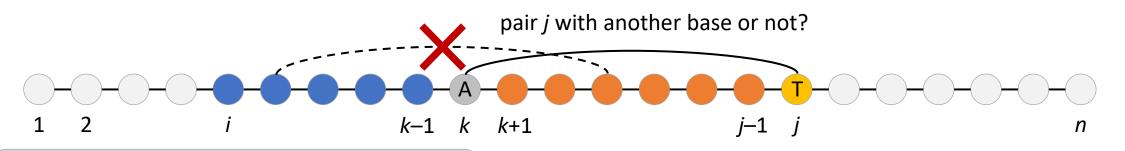**Recursive algorithm** (implement w/ dynamic programming)**:**
OPT($i,j$) = max of:
   OPT($i,j-1$),                *// don't form base pair with j*
   $\max_{i \leq k < j}$ 1 + OPT($i,k-1$) + OPT($k+1,j-1$) *// form k,j base pair*

12

# Computing maximally bound unpseudoknotted secondary structure in polynomial time

pair $j$ with another base or not?



**Recursive solution**:
- Strand length is $n$.
- For $1 \leq i \leq j \leq n$, let OPT$(i,j)$ = max base pairs possible using **only** bases $i$ through $j$.
- <u>Question</u>: do we pair base $j$ with some other base between $i$ and $j{-}1$?
- If *not*, recursively, the optimal value is:
  - OPT$(i,j)$ = OPT$(i,j{-}1)$
- If we pair $j$ with $k$, **nesting property** implies no base pair can form between any base in $[i,\ldots k{-}1]$ and any base in $[k{+}1,j{-}1]$
- Recursively, optimal value depends on:
  - OPT$(i,k{-}1)$ and OPT$(k{+}1,j{-}1)$

**Recursive algorithm** (implement w/ dynamic programming):
OPT$(i,j)$ = max of:

only if $k$ and $j$ are complementary bases

  OPT$(i,j{-}1)$,     // *don't form base pair with $j$*

  $\max_{i \leq k < j}$   $1$ + OPT$(i,k{-}1)$ + OPT$(k{+}1,j{-}1)$   // *form $k,j$ base pair*

12

# Computing maximally bound unpseudoknotted secondary structure in polynomial time



**Recursive solution**:
- Strand length is $n$.
- For $1 \leq i \leq j \leq n$, let OPT($i,j$) = max base pairs possible using **only** bases $i$ through $j$.
- <u>Question</u>: do we pair base $j$ with some other base between $i$ and $j{-}1$?
- If *not*, recursively, the optimal value is:
  - OPT($i,j$) = OPT($i,j{-}1$)
- If we pair $j$ with $k$, **nesting property** implies no base pair can form between any base in $[i,\dots k{-}1]$ and any base in $[k{+}1,j{-}1]$
- Recursively, optimal value depends on:
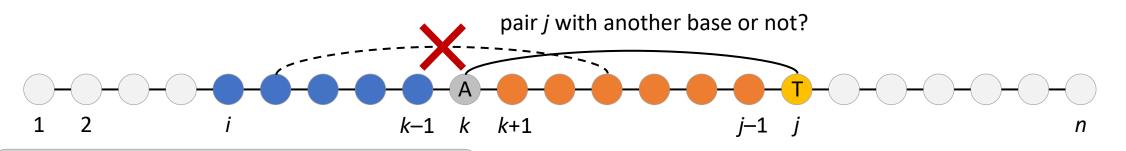  - OPT($i,k{-}1$) and OPT($k{+}1,j{-}1$)

**Recursive algorithm** (implement w/ dynamic programming):
OPT($i,j$) = max of:
   OPT($i,j{-}1$),     *// don't form base pair with j*
   $\max\limits_{i \leq k < j}$   1 + OPT($i,k{-}1$) + OPT($k{+}1,j{-}1$)  *// form k,j base pair*
base case: OPT($i,i$) = 0

only if $k$ and $j$ are complementary bases

12

# Computing maximally bound unpseudoknotted secondary structure in polynomial time

pair $j$ with another base or not?



1  2  ... $i$ ... $k-1$  $k$  $k+1$ ... $j-1$  $j$ ... $n$

**Recursive solution:**
- Strand length is $n$.
- For $1 \le i \le j \le n$, let OPT($i,j$) = max base pairs possible using **only** bases $i$ through $j$.
- <u>Question</u>: do we pair base $j$ with some other base between $i$ and $j-1$?
- If *not*, recursively, the optimal value is:
  - OPT($i,j$) = OPT($i,j-1$)
- If we pair $j$ with $k$, **nesting property** implies no base pair can form between any base in $[i,... k-1]$ and any base in $[k+1,j-1]$
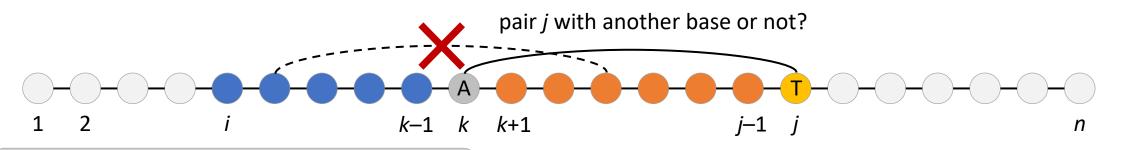- Recursively, optimal value depends on:
  - OPT($i,k-1$) and OPT($k+1,j-1$)

**Recursive algorithm** (implement w/ dynamic programming)**:**
OPT($i,j$) = max of:      only if $k$ and $j$ are complementary bases
   OPT($i,j-1$),             // *don't form base pair with j*
   $\max_{i \le k < j}$ 1 + OPT($i,k-1$) + OPT($k+1,j-1$)   // *form k,j base pair*
base case: OPT($i,i$) = 0
optimal value for whole strand = OPT($1,n$)

12

# Computing maximally bound unpseudoknotted secondary structure in polynomial time

pair $j$ with another base or not?



1  2  $i$  $k–1$  $k$  $k+1$  $j–1$  $j$  $n$

**Recursive solution**:
- Strand length is $n$.
- For $1 \leq i \leq j \leq n$, let OPT($i,j$) = max base pairs possible using **only** bases $i$ through $j$.
- <u>Question</u>: do we pair base $j$ with some other base between $i$ and $j–1$?
- If *not*, recursively, the optimal value is:
  - OPT($i,j$) = OPT($i,j–1$)
- If we pair $j$ with $k$, **nesting property** implies no base pair can form between any base in $[i,… k–1]$ and any base in $[k+1,j–1]$
- Recursively, optimal value depends on:
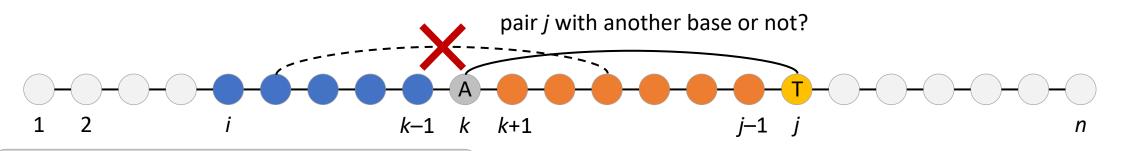  - OPT($i,k–1$) and OPT($k+1,j–1$)

**Recursive algorithm** (implement w/ dynamic programming):
OPT($i,j$) = max of: ⎫ only if $k$ and $j$ are complementary bases
   OPT($i,j–1$),          *// don't form base pair with j*
   $\max_{i \leq k < j}$  1 + OPT($i,k–1$) + OPT($k+1,j–1$)  *// form k,j base pair*
base case: OPT($i,i$) = 0
optimal value for whole strand = OPT($1,n$)

**Running time**:
There are $O(n^2)$ subproblems: choices $i,j$ with $1 \leq i < j \leq n$.

12

# Computing maximally bound unpseudoknotted secondary structure in polynomial time

pair $j$ with another base or not?



1  2        $i$           $k$–1  $k$  $k$+1            $j$–1  $j$                    $n$

**Recursive solution**:
- Strand length is $n$.
- For $1 \le i \le j \le n$, let OPT($i,j$) = max base pairs possible using **only** bases $i$ through $j$.
- <u>Question</u>: do we pair base $j$ with some other base between $i$ and $j$–1?
- If *not*, recursively, the optimal value is:
  - OPT($i,j$) = OPT($i,j$–1)
- If we pair $j$ with $k$, **nesting property** implies no base pair can form between any base in [$i,\dots k$–1] and any base in [$k$+1,$j$–1]
- Recursively, optimal value depends on:
  - OPT($i,k$–1) and OPT($k$+1,$j$–1)
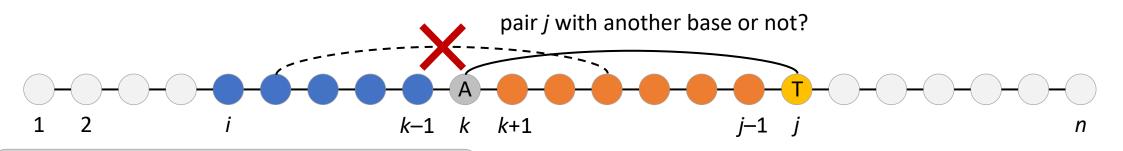
**Recursive algorithm** (implement w/ dynamic programming):
OPT($i,j$) = max of:                    only if $k$ and $j$ are complementary bases
   OPT($i,j$–1),                            *// don't form base pair with $j$*
   $\max_{i \le k < j}$  1 + OPT($i,k$–1) + OPT($k$+1,$j$–1)  *// form $k,j$ base pair*
base case: OPT($i,i$) = 0
optimal value for whole strand = OPT($1,n$)

**Running time**:
There are $O(n^2)$ subproblems: choices $i,j$ with $1 \le i < j \le n$.
Each takes time $O(n)$ to search all values of $k$, so $O(n^3)$ total.

# Computing maximally bound unpseudoknotted secondary structure in polynomial time

This gives optimal *value*: how to find <u>actual secondary structure</u>?

pair *j* with another base or not?



1  2        *i*              *k*−1  *k*  *k*+1              *j*−1  *j*              *n*

**Recursive solution**:
- Strand length is *n*.
- For 1 ≤ *i* ≤ *j* ≤ *n*, let OPT(*i*,*j*) = max base pairs possible using **only** bases *i* through *j*.
- <u>Question</u>: do we pair base *j* with some other base between *i* and *j*−1?
- If *not*, recursively, the optimal value is:
  - OPT(*i*,*j*) = OPT(*i*,*j*−1)
- If we pair *j* with *k*, **nesting property** implies no base pair can form between any base in [*i*,… *k*−1] and any base in [*k*+1,*j*−1]
- Recursively, optimal value depends on:
  - OPT(*i*,*k*−1) and OPT(*k*+1,*j*−1)

**Recursive algorithm** (implement w/ dynamic programming):
OPT(*i*,*j*) = max of:
   OPT(*i*,*j*−1),               *// don't form base pair with j*
   max$_{i \le k < j}$  1 + OPT(*i*,*k*−1) + OPT(*k*+1,*j*−1)  *// form k,j base pair*
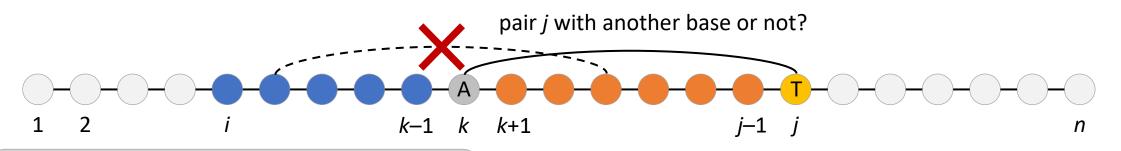base case: OPT(*i*,*i*) = 0
optimal value for whole strand = OPT(1,*n*)

only if *k* and *j* are complementary bases

**Running time**:
There are $O(n^2)$ subproblems: choices *i*,*j* with 1 ≤ *i* < *j* ≤ *n*.
Each takes time $O(n)$ to search all values of *k*, so $O(n^3)$ total.

# Example of dynamic programming algorithm

strand sequence =

# ATTGATC

# Example of dynamic programming algorithm

|       | A | T | T | G | A | T | C |
|-------|---|---|---|---|---|---|---|
| *i/j* | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| A 1   |   |   |   |   |   |   |   |
| T 2   |   |   |   |   |   |   |   |
| T 3   |   |   |   |   |   |   |   |
| G 4   |   |   |   |   |   |   |   |
| A 5   |   |   |   |   |   |   |   |
| T 6   |   |   |   |   |   |   |   |
| C 7   |   |   |   |   |   |   |   |

strand sequence =

# ATTGATC

# Example of dynamic programming algorithm

|       |   | A | T | T | G | A | T | C |
|-------|---|---|---|---|---|---|---|---|
| *i/j* |   | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| A | 1 |   |   |   |   |   |   |   |
| T | 2 | ✗ |   |   |   |   |   |   |
| T | 3 | ✗ | ✗ |   |   |   |   |   |
| G | 4 | ✗ | ✗ | ✗ |   |   |   |   |
| A | 5 | ✗ | ✗ | ✗ | ✗ |   |   |   |
| T | 6 | ✗ | ✗ | ✗ | ✗ | ✗ |   |   |
| C | 7 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |   |

strand sequence =

# ATTGATC

# Example of dynamic programming algorithm

|  | | A | T | T | G | A | T | C |
|---|---|---|---|---|---|---|---|---|
| | *i/j* | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| A | 1 | | | | | | | |
| T | 2 | ✗ | | | | | | |
| T | 3 | ✗ | ✗ | | | | | |
| G | 4 | ✗ | ✗ | ✗ | | | | |
| A | 5 | ✗ | ✗ | ✗ | ✗ | | | |
| T | 6 | ✗ | ✗ | ✗ | ✗ | ✗ | | |
| C | 7 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | |

strand sequence =

# ATTGATC

base cases

# Example of dynamic programming algorithm

|       | A | T | T | G | A | T | C |
|-------|---|---|---|---|---|---|---|
| *i/j* | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| **A** 1 | 0 |   |   |   |   |   |   |
| **T** 2 | ✗ |   |   |   |   |   |   |
| **T** 3 | ✗ | ✗ |   |   |   |   |   |
| **G** 4 | ✗ | ✗ | ✗ |   |   |   |   |
| **A** 5 | ✗ | ✗ | ✗ | ✗ |   |   |   |
| **T** 6 | ✗ | ✗ | ✗ | ✗ | ✗ |   |   |
| **C** 7 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |   |

strand sequence =

# ATTGATC

base cases

# Example of dynamic programming algorithm

|  | i/j | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|-----|---|---|---|---|---|---|---|
|  |  | A | T | T | G | A | T | C |
| A | 1 | 0 |  |  |  |  |  |  |
| T | 2 | ✗ | 0 |  |  |  |  |  |
| T | 3 | ✗ | ✗ |  |  |  |  |  |
| G | 4 | ✗ | ✗ | ✗ |  |  |  |  |
| A | 5 | ✗ | ✗ | ✗ | ✗ |  |  |  |
| T | 6 | ✗ | ✗ | ✗ | ✗ | ✗ |  |  |
| C | 7 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |  |

strand sequence =

# ATTGATC

base cases

13

# Example of dynamic programming algorithm



strand sequence =

# ATTGATC

base cases

# Example of dynamic programming algorithm

|  | i/j | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|
| | | A | T | T | G | A | T | C |
| A | 1 | 0 | | | | | | |
| T | 2 | ✕ | 0 | | | | | |
| T | 3 | ✕ | ✕ | 0 | | | | |
| G | 4 | ✕ | ✕ | ✕ | 0 | | | |
| A | 5 | ✕ | ✕ | ✕ | ✕ | | | |
| T | 6 | ✕ | ✕ | ✕ | ✕ | ✕ | | |
| C | 7 | ✕ | ✕ | ✕ | ✕ | ✕ | ✕ | |

strand sequence =

# ATTGATC

base cases

# Example of dynamic programming algorithm

|  | i/j | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|
| | | A | T | T | G | A | T | C |
| A | 1 | 0 | | | | | | |
| T | 2 | ✕ | 0 | | | | | |
| T | 3 | ✕ | ✕ | 0 | | | | |
| G | 4 | ✕ | ✕ | ✕ | 0 | | | |
| A | 5 | ✕ | ✕ | ✕ | ✕ | 0 | | |
| T | 6 | ✕ | ✕ | ✕ | ✕ | ✕ | | |
| C | 7 | ✕ | ✕ | ✕ | ✕ | ✕ | ✕ | |

strand sequence =

# ATTGATC

base cases

# Example of dynamic programming algorithm

|  | i/j | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|
|  |  | A | T | T | G | A | T | C |
| A | 1 | 0 |  |  |  |  |  |  |
| T | 2 | ✗ | 0 |  |  |  |  |  |
| T | 3 | ✗ | ✗ | 0 |  |  |  |  |
| G | 4 | ✗ | ✗ | ✗ | 0 |  |  |  |
| A | 5 | ✗ | ✗ | ✗ | ✗ | 0 |  |  |
| T | 6 | ✗ | ✗ | ✗ | ✗ | ✗ | 0 |  |
| C | 7 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |  |

strand sequence =

## ATTGATC

base cases

# Example of dynamic programming algorithm

| i/j | A 1 | T 2 | T 3 | G 4 | A 5 | T 6 | C 7 |
|-----|-----|-----|-----|-----|-----|-----|-----|
| A 1 | 0 | | | | | | |
| T 2 | ✕ | 0 | | | | | |
| T 3 | ✕ | ✕ | 0 | | | | |
| G 4 | ✕ | ✕ | ✕ | 0 | | | |
| A 5 | ✕ | ✕ | ✕ | ✕ | 0 | | |
| T 6 | ✕ | ✕ | ✕ | ✕ | ✕ | 0 | |
| C 7 | ✕ | ✕ | ✕ | ✕ | ✕ | ✕ | 0 |

strand sequence =

# ATTGATC

base cases

13

# Example of dynamic programming algorithm



strand sequence =

## ATTGATC

base cases

recursive cases with complementary bases

# Example of dynamic programming algorithm



strand sequence =

# ATTGATC

base cases

recursive cases with complementary bases

# Example of dynamic programming algorithm

strand sequence =

# ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

13

# Example of dynamic programming algorithm

strand sequence =

## ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

# Example of dynamic programming algorithm

strand sequence =

## ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

13

# Example of dynamic programming algorithm



strand sequence =

# ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

# Example of dynamic programming algorithm

|  | i/j | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|-----|---|---|---|---|---|---|---|
|  |  | A | T | T | G | A | T | C |
| A | 1 | 0 | 1 |  |  |  |  |  |
| T | 2 | ✕ | 0 | 0 |  |  |  |  |
| T | 3 | ✕ | ✕ | 0 | 0 |  |  |  |
| G | 4 | ✕ | ✕ | ✕ | 0 | 0 |  |  |
| A | 5 | ✕ | ✕ | ✕ | ✕ | 0 |  |  |
| T | 6 | ✕ | ✕ | ✕ | ✕ | ✕ | 0 |  |
| C | 7 | ✕ | ✕ | ✕ | ✕ | ✕ | ✕ | 0 |

strand sequence =

# ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

13

# Example of dynamic programming algorithm

|     | i/j | A | T | T | G | A | T | C |
|-----|-----|---|---|---|---|---|---|---|
|     |     | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| A   | 1   | 0 | 1 |   |   |   |   |   |
| T   | 2   | ✗ | 0 | 0 |   |   |   |   |
| T   | 3   | ✗ | ✗ | 0 | 0 |   |   |   |
| G   | 4   | ✗ | ✗ | ✗ | 0 | 0 |   |   |
| A   | 5   | ✗ | ✗ | ✗ | ✗ | 0 | 1 |   |
| T   | 6   | ✗ | ✗ | ✗ | ✗ | ✗ | 0 |   |
| C   | 7   | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | 0 |

strand sequence =

# ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

13

# Example of dynamic programming algorithm



strand sequence =

## ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

# Example of dynamic programming algorithm



|  | i/j | A 1 | T 2 | T 3 | G 4 | A 5 | T 6 | C 7 |
|---|---|---|---|---|---|---|---|---|
| A | 1 | 0 | 1 | ● | | | | |
| T | 2 | ✗ | 0 | 0 | | | | |
| T | 3 | ✗ | ✗ | 0 | 0 | | | |
| G | 4 | ✗ | ✗ | ✗ | 0 | 0 | | |
| A | 5 | ✗ | ✗ | ✗ | ✗ | 0 | 1 | |
| T | 6 | ✗ | ✗ | ✗ | ✗ | ✗ | 0 | 0 |
| C | 7 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | 0 |

strand sequence =

# ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

13

# Example of dynamic programming algorithm

strand sequence =

**ATTGATC**

base cases

recursive cases with complementary bases

recursive cases without complementary bases

Example of dynamic programming algorithm

strand sequence =

# ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

13

# Example of dynamic programming algorithm



strand sequence =
# ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

# Example of dynamic programming algorithm



strand sequence =

## ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

13

# Example of dynamic programming algorithm

strand sequence =

**ATTGATC**

base cases

recursive cases with complementary bases

recursive cases without complementary bases

13

# Example of dynamic programming algorithm

strand sequence =

**ATTGATC**

base cases

recursive cases with complementary bases

recursive cases without complementary bases

13

# Example of dynamic programming algorithm

strand sequence =

## ATTGATC

- base cases
- recursive cases with complementary bases
- recursive cases without complementary bases

# Example of dynamic programming algorithm



strand sequence =

## ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

Example of dynamic programming algorithm

strand sequence =

ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

13

# Example of dynamic programming algorithm

strand sequence =

# ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

# Example of dynamic programming algorithm

strand sequence =

# ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

Example of dynamic programming algorithm

strand sequence =

# ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

13

# Example of dynamic programming algorithm

strand sequence =

# ATTGATC

base cases

recursive cases with complementary bases

recursive cases without complementary bases

# Extensions to more realistic energy models

- base pairs on one strand must be separated by at least 4 other bases

# Extensions to more realistic energy models

- base pairs on one strand must be separated by at least 4 other bases
  - base case switches from OPT($i,i$) = 0 to OPT($i,j$)=0 if $j–i \leq 4$

# Extensions to more realistic energy models

- base pairs on one strand must be separated by at least 4 other bases
  - base case switches from OPT($i,i$) = 0 to OPT($i,j$)=0 if $j–i \leq 4$
- G/C twice as strong as A/T?

# Extensions to more realistic energy models

- base pairs on one strand must be separated by at least 4 other bases
  - base case switches from OPT($i$,$i$) = 0 to OPT($i$,$j$)=0 if $j$–$i$ ≤ 4
- G/C twice as strong as A/T?
  - $\max_{i \leq k < j}$ <span style="color:red">(1 if $k$,$j$ is A/T base pair, else 2)</span> + OPT($i$,$k$–1) + OPT($k$+1,$j$–1)

# Extensions to more realistic energy models

- base pairs on one strand must be separated by at least 4 other bases
  - base case switches from OPT($i$,$i$) = 0 to OPT($i$,$j$)=0 if $j$–$i$ ≤ 4
- G/C twice as strong as A/T?
  - $\max_{i \leq k < j}$ (1 if $k$,$j$ is A/T base pair, else 2) + OPT($i$,$k$–1) + OPT($k$+1,$j$–1)
- nearest-neighbor interaction?
  - $\max_{i \leq k < j}$ (more complex lookup here) + OPT($i$,$k$–1) + OPT($k$+1,$j$–1)

# Extensions to more realistic energy models

- base pairs on one strand must be separated by at least 4 other bases
  - base case switches from OPT($i,i$) = 0 to OPT($i,j$)=0 if $j–i ≤ 4$
- G/C twice as strong as A/T?
  - $\max_{i≤k<j}$ (1 if $k,j$ is A/T base pair, else 2) + OPT($i,k–1$) + OPT($k+1,j–1$)
- nearest-neighbor interaction?
  - $\max_{i≤k<j}$ (more complex lookup here) + OPT($i,k–1$) + OPT($k+1,j–1$)
- multiple strands?
  - a $\Delta G_{assoc}$ term for each strand beyond the first one
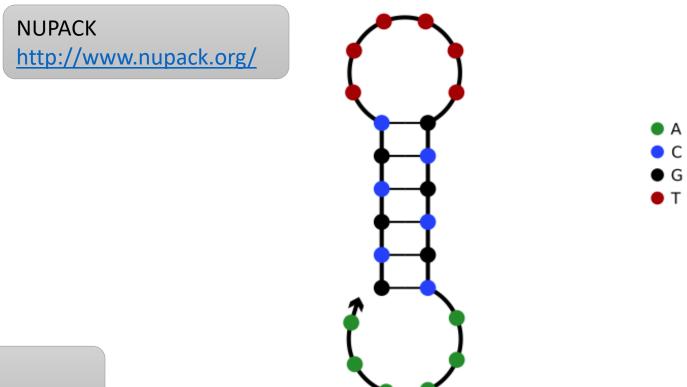
# Extensions to more realistic energy models

- base pairs on one strand must be separated by at least 4 other bases
  - base case switches from OPT($i$,$i$) = 0 to OPT($i$,$j$)=0 if $j–i \leq 4$
- G/C twice as strong as A/T?
  - $\max_{i \leq k < j}$ <span style="color:red">(1 if $k$,$j$ is A/T base pair, else 2)</span> + OPT($i$,$k–1$) + OPT($k+1$,$j–1$)
- nearest-neighbor interaction?
  - $\max_{i \leq k < j}$ <span style="color:red">(more complex lookup here)</span> + OPT($i$,$k–1$) + OPT($k+1$,$j–1$)
- multiple strands?
  - a $\Delta G_{assoc}$ term for each strand beyond the first one
- https://piercelab-caltech.github.io/nupack-docs/definitions/

# Software to compute minimum free energy DNA structures

NUPACK
http://www.nupack.org/

ViennaRNA
https://www.tbi.univie.ac.at/RNA/

MFE structure at 37.0 C

A
C
G
T

Free energy of secondary structure: -8.78 kcal/mol

# What is "free energy"?

A way to express <u>probability</u> of seeing a structure, in units of energy (kcal/mol).

Energy and probability are *exponentially* related.

# What is "free energy"?

A way to express <u>probability</u> of seeing a structure, in units of energy (kcal/mol).

Energy and probability are *exponentially* related.

- If *S* is a secondary structure, let Pr[*S*] denote probability of seeing it ("*at equilibrium*").

# What is "free energy"?

A way to express <u>probability</u> of seeing a structure, in units of energy (kcal/mol).

Energy and probability are *exponentially* related.

- If $S$ is a secondary structure, let Pr[$S$] denote probability of seeing it ("*at equilibrium*").
- At fixed temperature, ln(Pr[$S$]) ≈ Δ$G$($S$)   (*recall free energy ΔG(S) is negative*)

# What is "free energy"?

A way to express <u>probability</u> of seeing a structure, in units of energy (kcal/mol).

Energy and probability are *exponentially* related.

- If $S$ is a secondary structure, let Pr[$S$] denote probability of seeing it ("*at equilibrium*").
- At fixed temperature, ln(Pr[$S$]) ≈ $\Delta G(S)$   (*recall free energy $\Delta G(S)$ is negative*)
- Some constants: ln(Pr[$S$]) ≈ $\Delta G(S)/(RT)$, usually expressed as Pr[$S$] ∝ $e^{-\Delta G(S)/(RT)}$
  - $T$ = temperature in K (Kelvin), $R$ = Boltzmann's constant ≈ 0.001987204 kcal/mol/K

# What is "free energy"?

A way to express <u>probability</u> of seeing a structure, in units of energy (kcal/mol).

Energy and probability are *exponentially* related.

- If $S$ is a secondary structure, let Pr[$S$] denote probability of seeing it ("*at equilibrium*").
- At fixed temperature, ln(Pr[$S$]) ≈ $\Delta G(S)$   (*recall free energy $\Delta G(S)$ is negative*)
- Some constants: ln(Pr[$S$]) ≈ $\Delta G(S)/(RT)$, usually expressed as Pr[$S$] ∝ $e^{-\Delta G(S)/(RT)}$
    $T$ = temperature in K (Kelvin), $R$ = Boltzmann's constant ≈ 0.001987204 kcal/mol/K
- To convert $e^{-\Delta G(S)/(RT)}$ to a <u>probability</u>, need to normalize so they <u>sum to 1</u>.

# What is "free energy"?

A way to express <u>probability</u> of seeing a structure, in units of energy (kcal/mol).

Energy and probability are *exponentially* related.

- If $S$ is a secondary structure, let Pr[$S$] denote probability of seeing it ("*at equilibrium*").
- At fixed temperature, ln(Pr[$S$]) ≈ $\Delta G(S)$  (*recall free energy $\Delta G(S)$ is negative*)
- Some constants: ln(Pr[$S$]) ≈ $\Delta G(S)/(RT)$, usually expressed as Pr[$S$] ∝ $e^{-\Delta G(S)/(RT)}$
    $T$ = temperature in K (Kelvin), $R$ = Boltzmann's constant ≈ 0.001987204 kcal/mol/K
- To convert $e^{-\Delta G(S)/(RT)}$ to a <u>probability</u>, need to normalize so they <u>sum to 1</u>.
- For a DNA strand/set of DNA strands, let Ω denote set of all secondary structures.

# What is "free energy"?

A way to express <u>probability</u> of seeing a structure, in units of energy (kcal/mol).

Energy and probability are *exponentially* related.

- If $S$ is a secondary structure, let Pr[$S$] denote probability of seeing it ("*at equilibrium*").
- At fixed temperature, ln(Pr[$S$]) ≈ $\Delta G(S)$   (*recall free energy $\Delta G(S)$ is negative*)
- Some constants: ln(Pr[$S$]) ≈ $\Delta G(S)/(RT)$, usually expressed as Pr[$S$] ∝ $e^{-\Delta G(S)/(RT)}$
     $T$ = temperature in K (Kelvin), $R$ = Boltzmann's constant ≈ 0.001987204 kcal/mol/K
- To convert $e^{-\Delta G(S)/(RT)}$ to a <u>probability</u>, need to normalize so they <u>sum to 1</u>.
- For a DNA strand/set of DNA strands, let Ω denote set of all secondary structures.

**Definition**: The <u>partition function</u>
of Ω is $Q = \Sigma_{S \in \Omega} \, e^{-\Delta G(S)/(RT)}$.

# What is "free energy"?

A way to express <u>probability</u> of seeing a structure, in units of energy (kcal/mol).

Energy and probability are *exponentially* related.

- If $S$ is a secondary structure, let Pr[$S$] denote probability of seeing it ("*at equilibrium*").
- At fixed temperature, ln(Pr[$S$]) ≈ $\Delta G(S)$   (*recall free energy $\Delta G(S)$ is negative*)
- Some constants: ln(Pr[$S$]) ≈ $\Delta G(S)/(RT)$, usually expressed as Pr[$S$] ∝ $e^{-\Delta G(S)/(RT)}$
    $T$ = temperature in K (Kelvin), $R$ = Boltzmann's constant ≈ 0.001987204 kcal/mol/K
- To convert $e^{-\Delta G(S)/(RT)}$ to a <u>probability</u>, need to normalize so they <u>sum to 1</u>.
- For a DNA strand/set of DNA strands, let Ω denote set of all secondary structures.

**Definition**: The <u>partition function</u> of Ω is $Q = \Sigma_{S \in \Omega}\ e^{-\ \Delta G(S)/(RT)}$.

For any secondary structure $S$, Pr[$S$] = $(1/Q) \cdot e^{-\Delta G(S)/(RT)}$.

# Minimum free energy versus complex free energy

**Recall**: For any secondary structure $S$, $\Pr[S] = (1/Q) \cdot e^{-\Delta G(S)/(RT)}$

Minimum free energy structure $S$ is the most likely structure.

# Minimum free energy versus complex free energy

**Recall**: For any secondary structure $S$, $Pr[S] = (1/Q) \cdot e^{-\Delta G(S)/(RT)}$

Minimum free energy structure $S$ is the most likely structure.

**Problem**: What if *most likely* structure $S$ is *not very likely*?



$Pr[\text{———————→}] = $ **0.2**, but

$Pr[\phantom{O}] = Pr[\phantom{O}] = Pr[\phantom{O}]$

$= Pr[\phantom{O}] = $ **0.199**

# Minimum free energy versus complex free energy

**Recall**: For any secondary structure $S$, $\Pr[S] = (1/Q) \cdot e^{-\Delta G(S)/(RT)}$

Minimum free energy structure $S$ is the most likely structure.

**Problem**: What if *most likely* structure $S$ is *not very likely*?



$\Pr[\longrightarrow] = \mathbf{0.2}$, but

$\Pr[\;] = \Pr[\;] = \Pr[\;]$

$= \Pr[\;] = \mathbf{0.199}$

This strand spends nearly 80% of its time bound.

# Minimum free energy versus complex free energy

**Recall**: For any secondary structure $S$, $\Pr[S] = (1/Q) \cdot e^{-\Delta G(S)/(RT)}$

<u>Minimum free energy</u> structure $S$ is the <u>most likely</u> structure.

**Problem**: What if *most likely* structure $S$ is *not very likely*?

**Solution**: Consider energy of *all secondary structures at once*.

$\Pr[$  $] = $ **0.2**, but

$\Pr[$  $] = \Pr[$  $] = \Pr[$  $]$

$= \Pr[$  $] = $ **0.199**   This strand spends nearly 80% of its time bound.

# Minimum free energy versus complex free energy

**Recall**: For any secondary structure $S$, $\Pr[S] = (1/Q) \cdot e^{-\Delta G(S)/(RT)}$

<u>Minimum free energy</u> structure $S$ is the <u>most likely</u> structure.

**Problem**: What if *most likely* structure $S$ is *not very likely*?

**Solution**: Consider energy of *all secondary structures at once*.

$\Pr[$  $] = $ **0.2**, but

$\Pr[$  $] = \Pr[$  $] = \Pr[$  $]$

$= \Pr[$  $] = $ **0.199**

This strand spends nearly 80% of its time bound.

**Definition**: The <u>complex free energy</u> of $\Omega$ is $\Delta G = -RT \ln Q$.
Intuitively captures how much we expect strand to be bound/structured: higher (closer to 0) means more unstructured.

https://piercelab-caltech.github.io/nupack-docs/definitions/#complex-free-energy

17

# Minimum free energy versus complex free energy

**Recall**: For any secondary structure $S$,
$$\Pr[S] = (1/Q) \cdot e^{-\Delta G(S)/(RT)}$$

Minimum free energy structure $S$ is the most likely structure.

**Problem**: What if *most likely* structure $S$ is *not very likely*?

**Solution**: Consider energy of *all secondary structures at once*.

$\Pr[$  $] = $ **0.2**, but

$\Pr[$  $] = \Pr[$  $] = \Pr[$  $]$

$= \Pr[$  $] = $ **0.199**

This strand spends nearly 80% of its time bound.

**Definition**: The complex free energy of $\Omega$ is $\Delta G = -RT \ln Q$.
Intuitively captures how much we expect strand to be bound/structured: higher (closer to 0) means more unstructured.

https://piercelab-caltech.github.io/nupack-docs/definitions/#complex-free-energy

$\Delta G$ can also be computed in time $O(n^3)$.

17

# Example: DNA sequence design for single-stranded tiles

Given many single-stranded tiles with four domains each (lengths 10 and 11), assign DNA sequences to them satisfying:



Abbreviated list of constraints similar to those used in [*Diverse and robust molecular algorithms using reprogrammable DNA self-assembly*. Woods, Doty, Myhrvold, Hui, Zhou, Yin, Winfree. Nature 2019.]

# Example: DNA sequence design for single-stranded tiles

Given many single-stranded tiles with four domains each (lengths 10 and 11), assign DNA sequences to them satisfying:

- ∀strands $s$, $\Delta G(s) \geq -1.65$ kcal/mol

# Example: DNA sequence design for single-stranded tiles
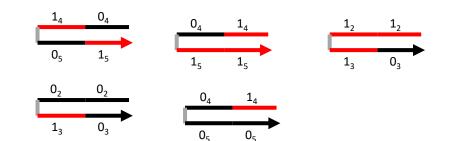
Given many single-stranded tiles with four domains each (lengths 10 and 11), assign DNA sequences to them satisfying:

- ∀strands $s$, $\Delta G(s) \geq -1.65$ kcal/mol

- ∀strand pairs $s,t$, $\Delta G(s,t) \geq -5.4$ kcal/mol if no complementary domains, $\geq -7.4$ kcal/mol otherwise



Abbreviated list of constraints similar to those used in [*Diverse and robust molecular algorithms using reprogrammable DNA self-assembly*. Woods, Doty, Myhrvold, Hui, Zhou, Yin, Winfree. Nature 2019.]
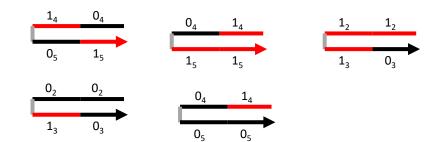
# Example: DNA sequence design for single-stranded tiles

Given many single-stranded tiles with four domains each (lengths 10 and 11), assign DNA sequences to them satisfying:

- ∀strands $s$, $\Delta G(s) \geq -1.65$ kcal/mol

- ∀strand pairs $s,t$, $\Delta G(s,t) \geq -5.4$ kcal/mol if no complementary domains, $\geq -7.4$ kcal/mol otherwise

- all domains end with A or T



Abbreviated list of constraints similar to those used in [*Diverse and robust molecular algorithms using reprogrammable DNA self-assembly*. Woods, Doty, Myhrvold, Hui, Zhou, Yin, Winfree. Nature 2019.]

# Example: DNA sequence design for single-stranded tiles

Given many single-stranded tiles with four domains each (lengths 10 and 11), assign DNA sequences to them satisfying:
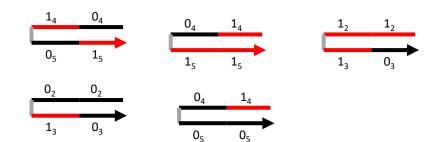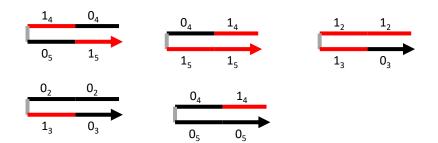
- $\forall$strands $s$, $\Delta G(s) \geq -1.65$ kcal/mol

- $\forall$strand pairs $s,t$, $\Delta G(s,t) \geq -5.4$ kcal/mol if no complementary domains, $\geq -7.4$ kcal/mol otherwise

- all domains end with A or T

- all domains have nearest-neighbor duplex energy between $-9.2$ and $-8.9$ kcal/mol



Abbreviated list of constraints similar to those used in [*Diverse and robust molecular algorithms using reprogrammable DNA self-assembly*. Woods, Doty, Myhrvold, Hui, Zhou, Yin, Winfree. Nature 2019.]

18

# Example: DNA sequence design for single-stranded tiles

Given many single-stranded tiles with four domains each (lengths 10 and 11), assign DNA sequences to them satisfying:
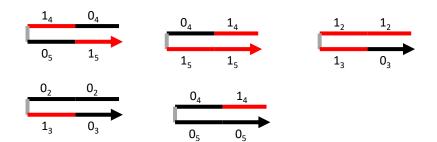


- $\forall$strands $s$, $\Delta G(s) \geq -1.65$ kcal/mol

- $\forall$strand pairs $s,t$, $\Delta G(s,t) \geq -5.4$ kcal/mol if no complementary domains, $\geq -7.4$ kcal/mol otherwise

- all domains end with A or T

- all domains have nearest-neighbor duplex energy between $-9.2$ and $-8.9$ kcal/mol

- tiles with even subscript domains on top have at most one G per domain (helps to satisfy first constraint)

Abbreviated list of constraints similar to those used in [*Diverse and robust molecular algorithms using reprogrammable DNA self-assembly*. Woods, Doty, Myhrvold, Hui, Zhou, Yin, Winfree. <u>Nature</u> 2019.]

# Example: DNA sequence design for single-stranded tiles

Given many single-stranded tiles with four domains each (lengths 10 and 11), assign DNA sequences to them satisfying:
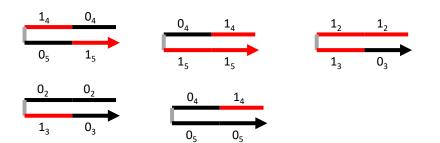
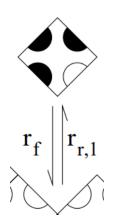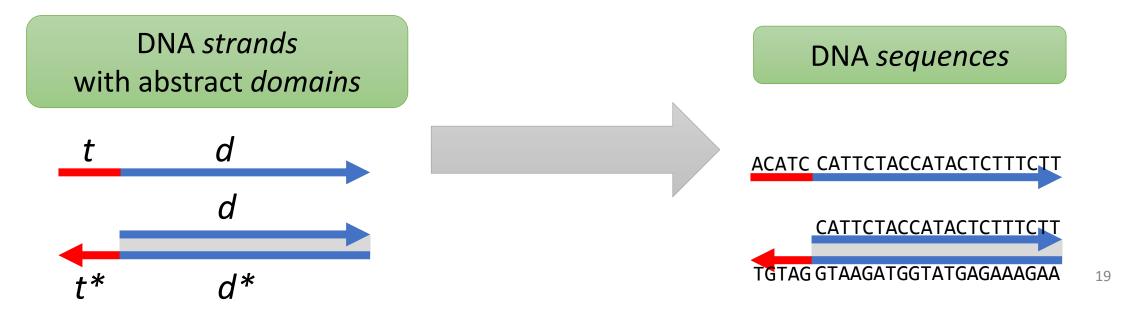- $\forall$strands $s$, $\Delta G(s) \geq -1.65$ kcal/mol

- $\forall$strand pairs $s,t$, $\Delta G(s,t) \geq -5.4$ kcal/mol if no complementary domains, $\geq -7.4$ kcal/mol otherwise

- all domains end with A or T

- all domains have nearest-neighbor duplex energy between $-9.2$ and $-8.9$ kcal/mol

- tiles with even subscript domains on top have at most one G per domain (helps to satisfy first constraint)

- pairs of domains $d_1,d_2$ that could result in one-domain mismatches during tile binding have $\Delta G(d_1,d_2) \geq -1.6$ kcal/mol

Abbreviated list of constraints similar to those used in [*Diverse and robust molecular algorithms using reprogrammable DNA self-assembly*. Woods, Doty, Myhrvold, Hui, Zhou, Yin, Winfree. Nature 2019.]

# DNA sequence design

- If we have DNA sequences, we can compute MFE/complex free energies of individual strands, pairs of strands, etc. in polynomial time.

- <u>DNA sequence design problem</u>: given abstract strands with abstract domains, assign concrete DNA sequences to the domains to satisfy a list of (experiment-specific) constraints.

- This is almost certainly **NP**-hard for any "reasonable" choice of constraints.



DNA *strands* with abstract *domains*

$t$   $d$

$d$

$t*$   $d*$

DNA *sequences*

ACATC CATTCTACCATACTCTTTCTT

CATTCTACCATACTCTTTCTT

TGTAG GTAAGATGGTATGAGAAAGAA

# Stochastic local search for finding DNA sequences

https://github.com/UC-Davis-molecular-computing/nuad

20

# Stochastic local search for finding DNA sequences

1. Assign DNA sequences randomly to domains.

https://github.com/UC-Davis-molecular-computing/nuad

# Stochastic local search for finding DNA sequences

1. Assign DNA sequences randomly to domains.
   - Each domain has a fixed length.

https://github.com/UC-Davis-molecular-computing/nuad

# Stochastic local search for finding DNA sequences

1. Assign DNA sequences randomly to domains.
   - Each domain has a fixed length.
   - Implicitly assign complement sequence to complement domains.

https://github.com/UC-Davis-molecular-computing/nuad

# Stochastic local search for finding DNA sequences

1. Assign DNA sequences randomly to domains.
   - Each domain has a fixed length.
   - Implicitly assign complement sequence to complement domains.
   - "Easy" single-domain constraints such as [*no GGGG*] or [*domains have A or T at each end*] can be automatically satisfied at this step.

https://github.com/UC-Davis-molecular-computing/nuad

# Stochastic local search for finding DNA sequences

1. Assign DNA sequences randomly to domains.
   - Each domain has a fixed length.
   - Implicitly assign complement sequence to complement domains.
   - "Easy" single-domain constraints such as [*no GGGG*] or [*domains have A or T at each end*] can be automatically satisfied at this step.
2. Check list of all constraints, tallying violations and "blaming" appropriate domains.

https://github.com/UC-Davis-molecular-computing/nuad

# Stochastic local search for finding DNA sequences

1. Assign DNA sequences randomly to domains.
   - Each domain has a fixed length.
   - Implicitly assign complement sequence to complement domains.
   - "Easy" single-domain constraints such as [*no GGGG*] or [*domains have A or T at each end*] can be automatically satisfied at this step.
2. Check list of all constraints, tallying violations and "blaming" appropriate domains.
   - For example, if a strand $s$ has too low $\Delta G(s)$, all domains on strand are blamed.

https://github.com/UC-Davis-molecular-computing/nuad

# Stochastic local search for finding DNA sequences

1. Assign DNA sequences randomly to domains.
   - Each domain has a fixed length.
   - Implicitly assign complement sequence to complement domains.
   - "Easy" single-domain constraints such as [*no GGGG*] or [*domains have A or T at each end*] can be automatically satisfied at this step.
2. Check list of all constraints, tallying violations and "blaming" appropriate domains.
   - For example, if a strand *s* has too low Δ$G(s)$, all domains on strand are blamed.
3. If no constraints violated, we're done!

https://github.com/UC-Davis-molecular-computing/nuad

# Stochastic local search for finding DNA sequences

1. Assign DNA sequences randomly to domains.
    - Each domain has a fixed length.
    - Implicitly assign complement sequence to complement domains.
    - "Easy" single-domain constraints such as [*no GGGG*] or [*domains have A or T at each end*] can be automatically satisfied at this step.
2. Check list of all constraints, tallying violations and "blaming" appropriate domains.
    - For example, if a strand $s$ has too low $\Delta G(s)$, all domains on strand are blamed.
3. If no constraints violated, we're done!
4. Otherwise, pick a domain $d$ at random in proportion to total "score" of violations it caused.

https://github.com/UC-Davis-molecular-computing/nuad

# Stochastic local search for finding DNA sequences

1. Assign DNA sequences randomly to domains.
    - Each domain has a fixed length.
    - Implicitly assign complement sequence to complement domains.
    - "Easy" single-domain constraints such as [*no GGGG*] or [*domains have A or T at each end*] can be automatically satisfied at this step.
2. Check list of all constraints, tallying violations and "blaming" appropriate domains.
    - For example, if a strand *s* has too low Δ*G*(*s*), all domains on strand are blamed.
3. If no constraints violated, we're done!
4. Otherwise, pick a domain *d* at random in proportion to total "score" of violations it caused.
5. Assign new random DNA sequence to *d*.

https://github.com/UC-Davis-molecular-computing/nuad

# Stochastic local search for finding DNA sequences

1. Assign DNA sequences randomly to domains.
   - Each domain has a fixed length.
   - Implicitly assign complement sequence to complement domains.
   - "Easy" single-domain constraints such as [*no GGGG*] or [*domains have A or T at each end*] can be automatically satisfied at this step.
2. Check list of all constraints, tallying violations and "blaming" appropriate domains.
   - For example, if a strand $s$ has too low $\Delta G(s)$, all domains on strand are blamed.
3. If no constraints violated, we're done!
4. Otherwise, pick a domain $d$ at random in proportion to total "score" of violations it caused.
5. Assign new random DNA sequence to $d$.
   - This change propagates through to all instances of $d$ and $d*$ on all strands.

https://github.com/UC-Davis-molecular-computing/nuad

# Stochastic local search for finding DNA sequences

1. Assign DNA sequences randomly to domains.
   - Each domain has a fixed length.
   - Implicitly assign complement sequence to complement domains.
   - "Easy" single-domain constraints such as [*no GGGG*] or [*domains have A or T at each end*] can be automatically satisfied at this step.
2. Check list of all constraints, tallying violations and "blaming" appropriate domains.
   - For example, if a strand *s* has too low Δ*G*(*s*), all domains on strand are blamed.
3. If no constraints violated, we're done!
4. Otherwise, pick a domain *d* at random in proportion to total "score" of violations it caused.
5. Assign new random DNA sequence to *d*.
   - This change propagates through to all instances of *d* and *d** on all strands.
6. Repeat step 2; if the new DNA sequence for *d* results in lower score of violations, keep it, otherwise, ignore it and pick a new random domain at step 4.

https://github.com/UC-Davis-molecular-computing/nuad

# Stochastic local search for finding DNA sequences

1. Assign DNA sequences randomly to domains.
   - Each domain has a fixed length.
   - Implicitly assign complement sequence to complement domains.
   - "Easy" single-domain constraints such as [*no GGGG*] or [*domains have A or T at each end*] can be automatically satisfied at this step.
2. Check list of all constraints, tallying violations and "blaming" appropriate domains.
   - For example, if a strand *s* has too low $\Delta G(s)$, all domains on strand are blamed.
3. If no constraints violated, we're done!
4. Otherwise, pick a domain *d* at random in proportion to total "score" of violations it caused.
5. Assign new random DNA sequence to *d*.
   - This change propagates through to all instances of *d* and *d\** on all strands.
6. Repeat step 2; if the new DNA sequence for *d* results in lower score of violations, keep it, otherwise, ignore it and pick a new random domain at step 4.
7. Repeat until no constraints are violated.

https://github.com/UC-Davis-molecular-computing/nuad

# Stochastic local search for finding DNA sequences

1. Assign DNA sequences randomly to domains.
   - Each domain has a fixed length.
   - Implicitly assign complement sequence to complement domains.
   - "Easy" single-domain constraints such as [*no GGGG*] or [*domains have A or T at each end*] can be automatically satisfied at this step.
2. Check list of all constraints, tallying violations and "blaming" appropriate domains.
   - For example, if a strand *s* has too low $\Delta G(s)$, all domains on strand are blamed.
3. If no constraints violated, we're done!
4. Otherwise, pick a domain *d* at random in proportion to total "score" of violations it caused.
5. Assign new random DNA sequence to *d*.
   - This change propagates through to all instances of *d* and *d** on all strands.
6. Repeat step 2; if the new DNA sequence for *d* results in lower score of violations, keep it, otherwise, ignore it and pick a new random domain at step 4.
7. Repeat until no constraints are violated.

https://github.com/UC-Davis-molecular-computing/nuad

Slow and unclever, but it works for any set of constraints.

20