# Algorithms for RNA design

Sebastian Will

École Polytechnique - Institute Polytechnique de Paris

AlgoSB 2025

[Thanks to Hua-Ting Yao for a number of slides]
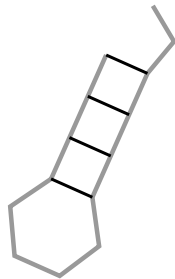
# RNA Design



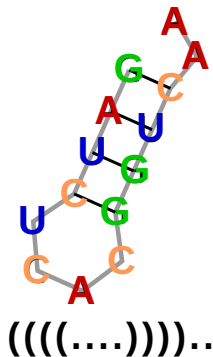RNA Design

GAUCUCACGGUCAA

A —— U
G —— C
Complementarity

(((( .... )))) ..

# RNA Design



GAUCUCACGGUCAA

Structure
prediction

Complementarity

(((( .... )))) ..

Structure design as "inverse folding"

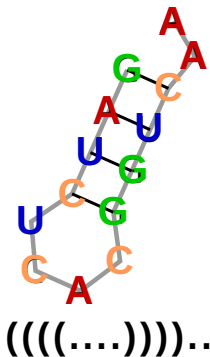# RNA Design



**GAUCUCACGGUCAA**

Structure prediction

A — U

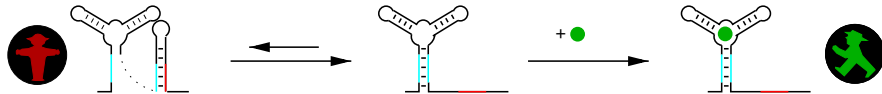G — C

Complementarity

((((....))))..

Structure design as "inverse folding"

More generally: Generating RNA sequences with desired functions

$\rightarrow$ biotechnology, RNA-based therapies, mRNA vaccines...
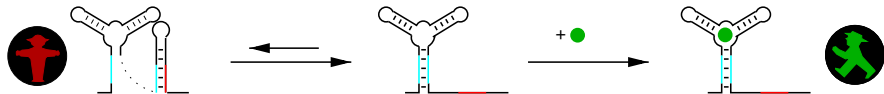
# Complex RNA design

**Design riboswitches for gene control**

# Complex RNA design

**Design riboswitches for gene control**



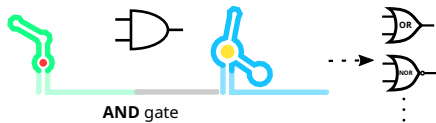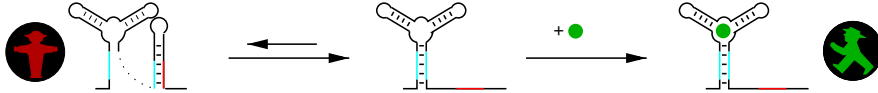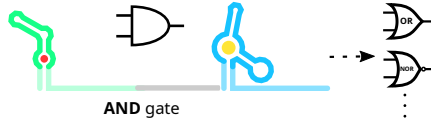Complex constructs: **AND-Riboswitch** (with G. Domin *et al.*, 2017)



**AND** gate

# Complex RNA design

## Design riboswitches for gene control



## Complex constructs: **AND-Riboswitch** (with G. Domin *et al.*, 2017)



**AND** gate

## Challenge: **Design for multiple target structures**



```
abcdefghijklmnopqrstuv
(((((.)).(((..))).))).
((.))((...))..(((..)))
....(((((..)))...)...
```

# Design of xrRNA riboswitch



PK 1

PK 2

Ring

xrRNA

PK 2

PK 1

Ring

degradation

with ligand

5'        CDS   3'

without ligand

5'        CDS   3'

[Leonhard Sidl, MT Wolfinger, HT Yao...]

# Design of SAM-I aptamer



- Aim for similarity with MSA
- Learn generative model from MSA (RBM)
- Compatibility with target structure (with PKs)
- Avoid off-targets



[Jorge Fernandez-de-Cossio-Diaz, 2024]

# Design of SAM-I riboswitch



[Jorge Fernandez-de-Cossio-Diaz, 2024]

Generation (Sampling):
- targets: compatibility, energy
- sequence similiarity

Refinement (Stochastic optimization):
- avoid off-targets
- relative stability bound/unbound

# RNA Design



GAUCUCACGGUCAA

RNA Design

A ══ U
G ══ C
Complementarity

((((....))))..

# RNA structure design: positive and negative

**Positive design:** Target a structure

$\rightarrow$ optimize affinity to target structures $t$

find sequence $\sigma$
with $E(\sigma, t) = \min_{\sigma'}(\sigma', t)$

extensions: multiple targets, properties, . . .

**Negative design**: Avoid all off-targets

$\rightarrow$ specificity for targets

find design $\sigma$,  s.t.
$E(\sigma, t) = \min_{t'} E(\sigma, t')$    $=: MFE(\sigma)$

refined objectives: probability, ensemble, . . .



good positive design
since all/most other sequences
energetically worse

**-18.8**

negative design
(by MFE criterion)

**-8.2**

but no negative design, since

**-23.8**

since all other structures
energetically worse, e.g.

**-7.8**    **-6.6**    ....

# Hello world of positive design!
## (single target structure, base pair energy)

**IN:** target structure $t$ of length $n$

"base pair energy" $E(\sigma, t) = \sum_{(i,j) \in t} e_{bp}(\sigma_i, \sigma_j) + \sum_{i \text{ unpaired in } t} e_u(\sigma_i)$

**Task:** sample from the Boltzmann distribution of sequences; i.e. sample $\sigma$ with probability

$$\Pr[\sigma] \sim \exp(-\beta E(\sigma, t)) \qquad \text{for some inverse temperature } \beta$$

# Hello world of positive design!
## (single target structure, base pair energy)

**IN:** target structure $t$ of length $n$

"base pair energy" $E(\sigma, t) = \sum_{(i,j) \in t} e_{bp}(\sigma_i, \sigma_j) + \sum_{i \text{ unpaired in } t} e_u(\sigma_i)$

**Task:** sample from the Boltzmann distribution of sequences; i.e. sample $\sigma$ with probability

$$\Pr[\sigma] \sim \exp(-\beta E(\sigma, t)) \qquad \text{for some inverse temperature } \beta$$

*Observation:* bases and base pairs can be generated *independently* of each other!

**Algorithm**

1) for each unpaired $i$: choose $\sigma_i$ with $Pr[\sigma_i] \sim exp(-\beta e_u(\sigma_i))$

2) for each pair $(i,j) \in t$: choose $\sigma_i$ and $\sigma_j$ with $Pr[\sigma_i, \sigma_j] \sim exp(-\beta e_{bp}(\sigma_i, \sigma_j))$

# Hello world of positive design!
## (single target structure, base pair energy)

**IN:** target structure $t$ of length $n$

"base pair energy" $E(\sigma, t) = \sum_{(i,j) \in t} e_{bp}(\sigma_i, \sigma_j) + \sum_{i \text{ unpaired in } t} e_u(\sigma_i)$

**Task:** sample from the Boltzmann distribution of sequences; i.e. sample $\sigma$ with probability

$$\Pr[\sigma] \sim \exp(-\beta E(\sigma, t)) \qquad \text{for some inverse temperature } \beta$$

*Observation:* bases and base pairs can be generated *independently* of each other!

**Algorithm**

1) for each unpaired $i$: choose $\sigma_i$ with $Pr[\sigma_i] \sim exp(-\beta e_u(\sigma_i))$

2) for each pair $(i, j) \in t$: choose $\sigma_i$ and $\sigma_j$ with $Pr[\sigma_i, \sigma_j] \sim exp(-\beta e_{bp}(\sigma_i, \sigma_j))$

🤔 What happens for complex energy models or multiple targets? (*dependencies!*)

# Infrared

General efficient framework for weighted constraint solving.
$\Rightarrow$ Rapid development of bioinformatics tools (including design)



Input

Functions: GC%, $E_1$, $E_2$, $E_3$

*Infrared solves a weighted form of constraint problems (CSP $\rightarrow$ Feature Networks)*
*It allows us to describe ("model") problems; then solves them automatically.*

[Hua-Ting Yao et al., 2024]

# Infrared

General efficient framework for weighted constraint solving.
$\Rightarrow$ Rapid development of bioinformatics tools (including design)



*Infrared solves a weighted form of constraint problems (CSP $\rightarrow$ Feature Networks)*
*It allows us to describe ("model") problems; then solves them automatically.*

[Hua-Ting Yao et al., 2024]

# Infrared

General efficient framework for weighted constraint solving.
$\Rightarrow$ Rapid development of bioinformatics tools (including design)
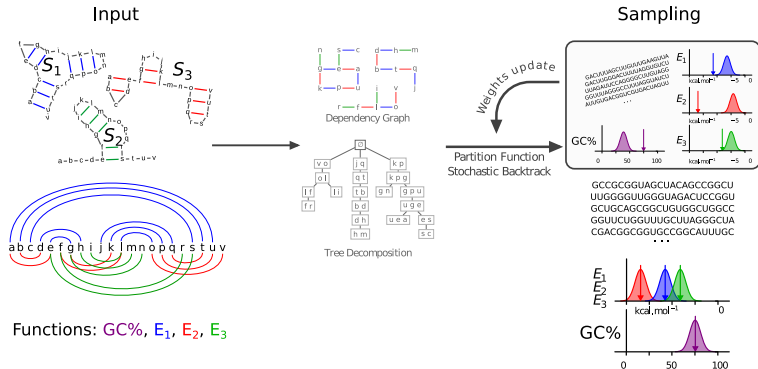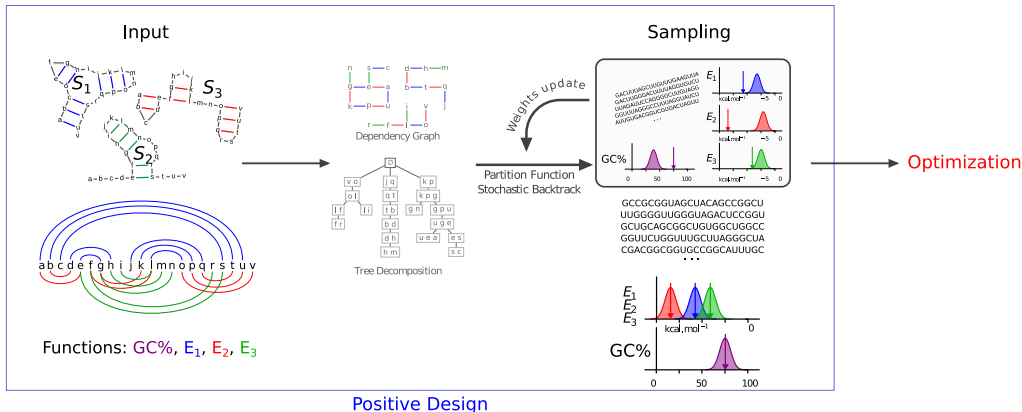


*Infrared solves a weighted form of constraint problems (CSP $\rightarrow$ Feature Networks)*
*It allows us to describe ("model") problems; then solves them automatically.*
[Hua-Ting Yao et al., 2024]

# Constraint Satisfaction Problems (CSPs)

Definition: A *CSP* is a tuple $(\mathcal{X}, \mathcal{D}, \mathcal{C})$, where

- $\mathcal{X} = \{X_1, \ldots, X_n\}$ is a set of *variables*
- $\mathcal{D} = \{D_1, \ldots, D_n\}$ is a set of corresponding finite *domains*
- $\mathcal{C}$ is a finite set of *constraints*

Each *constraint* $C$ is associated with $k$ variables.

**Solutions** of a CSP are **assignments** of domain values to the variables that satisfy all constraints (**valid assignments**).

# Constraint Satisfaction Problems (CSPs)

Definition: A *CSP* is a tuple $(\mathcal{X}, \mathcal{D}, \mathcal{C})$, where

- $\mathcal{X} = \{X_1, \ldots, X_n\}$ is a set of *variables*
- $\mathcal{D} = \{D_1, \ldots, D_n\}$ is a set of corresponding finite *domains*
- $\mathcal{C}$ is a finite set of *constraints*

Each *constraint* $C$ is associated with $k$ variables.

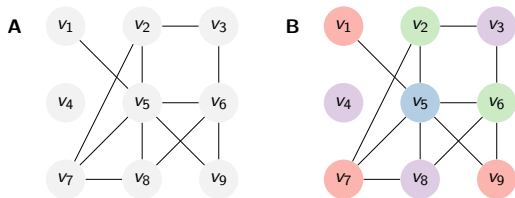**Solutions** of a CSP are **assignments** of domain values to the variables that satisfy all constraints (**valid assignments**).

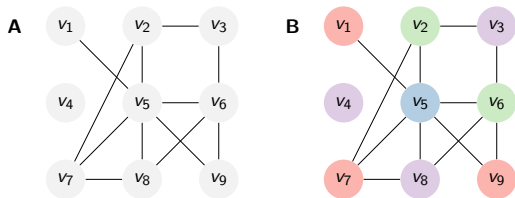General solving is NP-hard! Solving strategies?

- Generic, heuristic solving strategies: backtracking search + constraint propagation
- ... Infrared is specialized to problems with nearly tree-like dependencies

  we gain: efficient (fpt) exact optimization + controlled sampling!

# CSP Examples: Graph Coloring and N-Queens



**Constraints:** Adjacent nodes differ in color!

# CSP Examples: Graph Coloring and N-Queens



**Constraints:** Adjacent nodes differ in color!

$CSP = (\mathcal{X}, \mathcal{D}, \mathcal{C})$

- $\mathcal{X} = \{X_1, \ldots, X_9\}$

- $\mathcal{D} = \{X_1 \mapsto [1..4], \ldots, X_4 \mapsto [1..4]\}$

- $\mathcal{C} = \{X_i \neq X_j \mid i, j \text{ adjacent}\}$

```
model = Model(9,(1,4))
model.add_constraints(
    NotEquals(i,j) for i,j in edges)
```

# CSP Examples: Graph Coloring and N-Queens

**A**



**B**





$X_1 = 3$
$X_2 = 1$
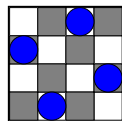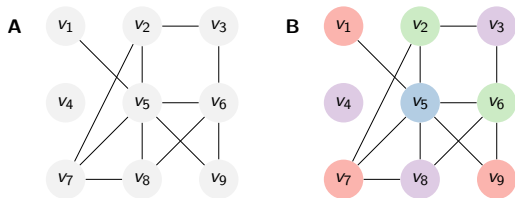$X_3 = 4$
$X_4 = 2$

**Constraints:** no attacks!



**Constraints:** Adjacent nodes differ in color!

$CSP = (\mathcal{X}, \mathcal{D}, \mathcal{C})$

- $\mathcal{X} = \{X_1, \ldots, X_9\}$

- $\mathcal{D} = \{X_1 \mapsto [1..4], \ldots, X_4 \mapsto [1..4]\}$

- $\mathcal{C} = \{X_i \neq X_j \mid i, j \text{adjacent}\}$

```
model = Model(9,(1,4))
model.add_constraints(
    NotEquals(i,j) for i,j in edges)
```

# CSP Examples: Graph Coloring and N-Queens

**A** 

**B** 



$X_1 = 3$
$X_2 = 1$
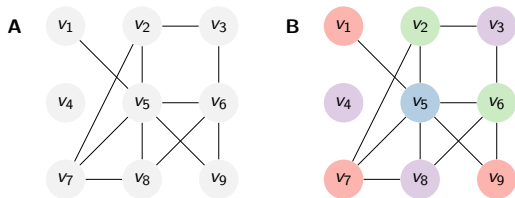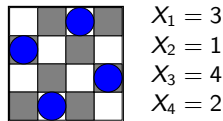$X_3 = 4$
$X_4 = 2$

**Constraints:** no attacks!



**Constraints:** Adjacent nodes differ in color!

$CSP = (\mathcal{X}, \mathcal{D}, \mathcal{C})$

- $\mathcal{X} = \{X_1, \ldots, X_9\}$

- $\mathcal{D} = \{X_1 \mapsto [1..4], \ldots, X_4 \mapsto [1..4]\}$

- $\mathcal{C} = \{X_i \neq X_j \mid i, j \text{ adjacent}\}$

```
model = Model(9,(1,4))
model.add_constraints(
    NotEquals(i,j) for i,j in edges)
```

$CSP = (\mathcal{X}, \mathcal{D}, \mathcal{C})$

- $\mathcal{X} = \{X_1, \ldots, X_4\}$

- $\mathcal{D} = \{X_1 \mapsto [1..4], \ldots, X_4 \mapsto [1..4]\}$

- $\mathcal{C} = \{X_i \neq X_j \mid 1 \leq i < j \leq 4\}$
  $\cup \{X_i + i \neq X_j + j \mid 1 \leq i < j \leq 4\}$
  $\cup \{X_i - i \neq X_j - j \mid 1 \leq i < j \leq 4\}$

```
model = Model(4,(1,4))
model.add_constraints(NotEquals(i,j)
    for i in range(4) for j in range(i+1,4))
model.add_constraint( ...
```

# Dependency graph and tree decomposition



```
model = Model(9,(1,4))
edges = [(1,5), (2,3), (2,5), ... ]
model.add_constraints(NotEquals(i,j)
    for i,j in edges)
```

# Dependency graph and tree decomposition



A

B

```
model = Model(9,(1,4))
edges = [(1,5), (2,3), (2,5), ... ]
model.add_constraints(NotEquals(i,j)
    for i,j in edges)
```

**Tree decomposition** $(T, \chi)$; $T = (V, E)$:

1. every variable occurs in one bag $v \in V$
2. for every constraint and function: there is one bag that contains its variables
3. for each variable: the bags containing it are connected

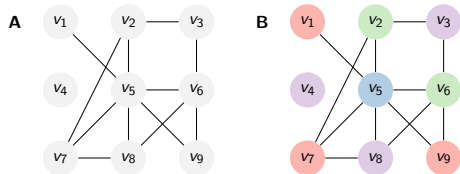**Tree width** = size of largest bag - 1

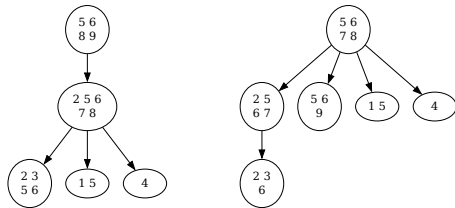**Remarks:**

# Dependency graph and tree decomposition



```
model = Model(9,(1,4))
edges = [(1,5), (2,3), (2,5), ... ]
model.add_constraints(NotEquals(i,j)
    for i,j in edges)
```
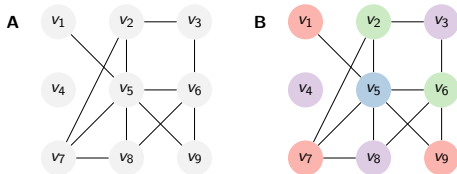


**Tree decomposition** $(T, \chi)$; $T = (V, E)$:

1. every variable occurs in one bag $v \in V$
2. for every constraint and function: there is one bag that contains its variables
3. for each variable: the bags containing it are connected

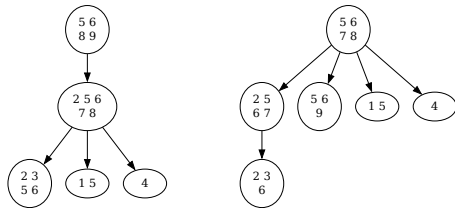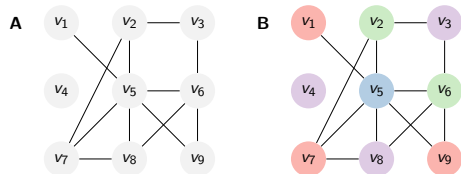**Tree width** = size of largest bag - 1

## Remarks:

- Conditions 1–3 allow solving by dynamic programming (solve from smaller to larger subtrees)
- Condition 2 → every constraint and function can be processed in some bag

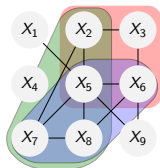# Adding functions → Objective function

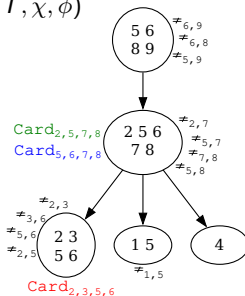*Make it more interesting by adding some functions*



A

$v_1$  $v_2$ — $v_3$

$v_4$  $v_5$ — $v_6$

$v_7$ — $v_8$  $v_9$

B

$v_1$  $v_2$  $v_3$

$v_4$  $v_5$  $v_6$

$v_7$  $v_8$  $v_9$

```
model = Model(9,(1,4))
model.add_constraints(NotEquals(i,j)
    for i,j in edges)

# extend by card feature
model.add_functions([Card(i,j,k,l)
    for i,j,k,l in fourcycles],'card')
```

**Cluster tree** $(T, \chi, \phi)$



$X_1$  $X_2$  $X_3$
$X_4$  $X_5$  $X_6$
$X_7$  $X_8$  $X_9$

$\begin{matrix} 5\ 6 \\ 8\ 9 \end{matrix}$  $\neq_{6,9}$ $\neq_{6,8}$ $\neq_{5,9}$

$\text{Card}_{2,5,7,8}$
$\text{Card}_{5,6,7,8}$  $\begin{matrix} 2\ 5\ 6 \\ 7\ 8 \end{matrix}$  $\neq_{2,7}$ $\neq_{5,7}$ $\neq_{7,8}$ $\neq_{5,8}$

$\neq_{2,3}$ $\neq_{3,6}$ $\neq_{5,6}$ $\neq_{2,5}$  $\begin{matrix} 2\ 3 \\ 5\ 6 \end{matrix}$  $\begin{matrix} 1\ 5 \end{matrix}$ $\neq_{1,5}$  $\begin{matrix} 4 \end{matrix}$

$\text{Card}_{2,3,5,6}$

# Feature Networks

*Feature networks add quality of solutions* → **features**.

Definition

A *Feature Network* is a tuple $\mathcal{N} = (\mathcal{X}, \mathcal{D}, \mathcal{C}, \mathcal{F})$, where

- $\mathcal{X} = \{X_1, \ldots, X_n\}$ is a set of *variables*
- $\mathcal{D} = \{D_1, \ldots, D_n\}$ is a set of corresponding finite *domains*
- $\mathcal{C}$ is a finite set of *constraints*
- $\mathcal{F}$ is a finite set of *features*, which consist of *feature functions*

*Features* . . .

- evaluate assignments as $F(x) = \sum_{f \in F} f(x)$
- define the *evaluation function* $E_{\mathcal{N}}(x, \alpha) = \sum_{F \in \mathcal{F}} \alpha_F F(x)$     for weights $\alpha_F$

# Infrared solves the sampling problem

Problem (Assignment sampling)

INPUT: *Feature Network $\mathcal{N}$, feature weights $\alpha$*
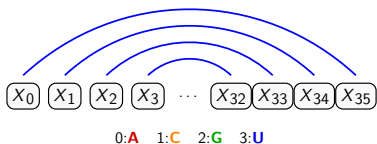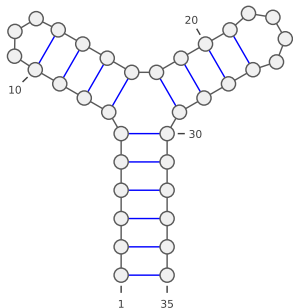
OUTPUT: *Valid assignment $x \in \mathcal{A}_\mathcal{X}$ generated with a probability that is proportional to its Boltzmann weight*

$$\mathbb{P}(x) \propto \exp\left(E_\mathcal{N}(x, \alpha)\right).$$

$\alpha = (\alpha_F)_{F \in \mathcal{F}}$ vector of weights

Evaluation function: $E_\mathcal{N}(x, \alpha) = \sum_{F \in \mathcal{F}} \alpha_F F(x)$.
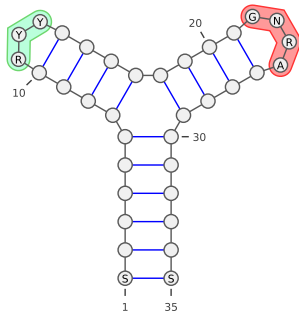
# Modeling: Single structure design



```python
import infrared as ir
import infrared.rna as rna

model = ir.Model(35, 4)   0:A  1:C  2:G  3:U

target = "(((((((((((...)))))(((((....)))))))))))"
model.add_constraints(rna.BPComp(i, j)   AU, CG, ...
    for (i, j) in rna.parse(target))
```



0:**A**  1:**C**  2:**G**  3:**U**

```python
sampler = ir.Sampler(model)
samples = [sampler.sample() for _ in range(1000)]
```

# Modeling: Single structure design



```python
import infrared as ir
import infrared.rna as rna

model = ir.Model(35, 4)   0:A  1:C  2:G  3:U

target = "(((((((((((...)))(((....))))))))))))"
model.add_constraints(rna.BPComp(i, j)   AU, CG, ···
    for (i, j) in rna.parse(target))

N:ACGU  S:CG  R:AG  Y:CU
iupac_seq = "SNNNNNNNNNRYYNNNNNNNNGNRANNNNNNNNNS"
for i, x in enumerate(iupac_seq):
    model.add_constraints(
        ir.ValueIn(i, rna.iupacvalues(x)))

sampler = ir.Sampler(model)
samples = [sampler.sample() for _ in range(1000)]
```

# Modeling: Single structure design



```python
import infrared as ir
import infrared.rna as rna

model = ir.Model(35, 4)  0:A 1:C 2:G 3:U

target = "((((((((((...))))(((....))))))))))"
model.add_constraints(rna.BPComp(i, j)   AU, CG, ···
    for (i, j) in rna.parse(target))


N:ACGU  S:CG  R:AG  Y:CU
iupac_seq = "SNNNNNNNNNRYYNNNNNNNNNGNRANNNNNNNNNS"
for i, x in enumerate(iupac_seq):
    model.add_constraints(
        ir.ValueIn(i, rna.iupacvalues(x)))


sampler = ir.Sampler(model)
samples = [sampler.sample() for _ in range(1000)]
```
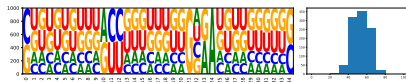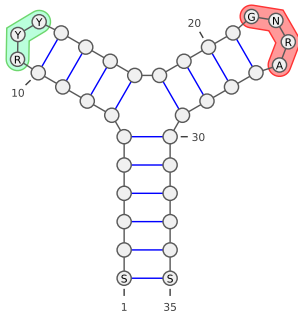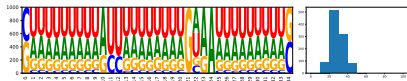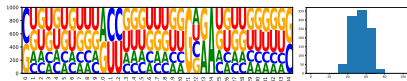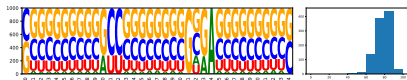
# Control GC-content



$\alpha = -1$



$\alpha = 0$



$\alpha = +1$

Method 1:

```
model.add_functions([rna.GCCont(i))    CG : 1   AU : 0
    for i in range(n)], 'gc')
model.set_feature_weight(α, 'gc')

sampler = ir.Sampler(model)
samples = [sampler.sample() for _ in range(1000)]
```

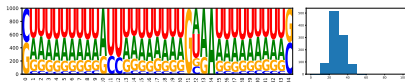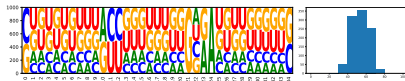# Control GC-content



$\alpha = -1$

$\alpha = 0$

$\alpha = +1$

Method 1:

```
model.add_functions([rna.GCCont(i))   CG : 1  AU : 0
    for i in range(n)], 'gc')
model.set_feature_weight(α, 'gc')

sampler = ir.Sampler(model)
samples = [sampler.sample() for _ in range(1000)]
```
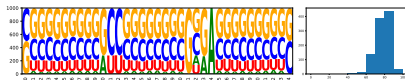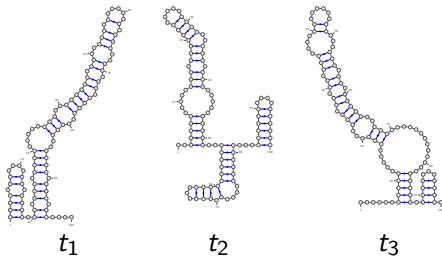
## Method 2 (Targeted sampling):

```
sampler = ir.Sampler(model)
sampler.set_target( 0.75 * n, 0.01 * n, 'gc' )
samples = [sampler.targeted_sample()
    for _ in range(1000)]     Automatically learn α
```
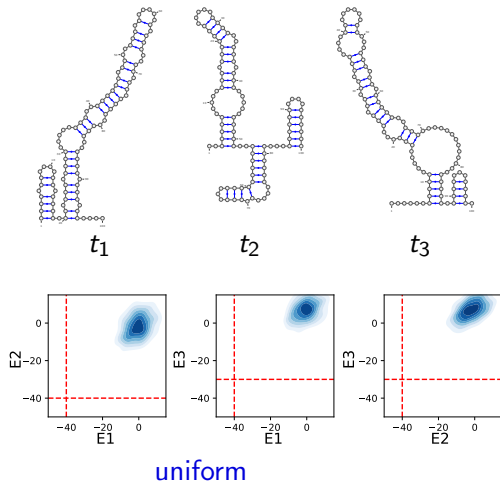
# Multitarget design



$t_1$          $t_2$          $t_3$

```python
model = ir.Model(n,4)

for k, target in enumerate(targets):
    bps = rna.parse(target)
    model.add_constraints(rna.BPComp(i, j)
        for (i, j) in bps)
```

# Multidimensional Boltzmann sampling



$t_1$     $t_2$     $t_3$



uniform

```python
model = ir.Model(n,4)
for k, target in enumerate(targets):
    bps = rna.parse(target)
    model.add_constraints(rna.BPComp(i, j)
        for (i, j) in bps)
                                    Simplified energy model
    model.add_functions([rna.BPEnergy(i, j)
        for (i, j) in bps], f'energy{k}')


for k, target in enumerate(targets):
    model.add_feature(f'E{k+1}', f'energy{k}',
        lambda sample, target=target:
            energy_of_struct(sample, target))
                    ViennaRNA energy model

sampler = ir.Sampler(model)
```

# Multidimensional Boltzmann sampling



$t_1$      $t_2$      $t_3$

uniform targeted

```python
model = ir.Model(n,4)
for k, target in enumerate(targets):
    bps = rna.parse(target)
    model.add_constraints(rna.BPComp(i, j)
        for (i, j) in bps)
                                    Simplified energy model
    model.add_functions([rna.BPEnergy(i, j)
        for (i, j) in bps], f'energy{k}')


for k, target in enumerate(targets):
    model.add_feature(f'E{k+1}', f'energy{k}',
        lambda sample, target=target:
            energy_of_struct(sample, target))
                      ViennaRNA energy model

sampler = ir.Sampler(model)
sampler.set_target( -40, 0.5, 'E1')
sampler.set_target( -40, 0.5, 'E2')
sampler.set_target( -30, 0.5, 'E3')
```
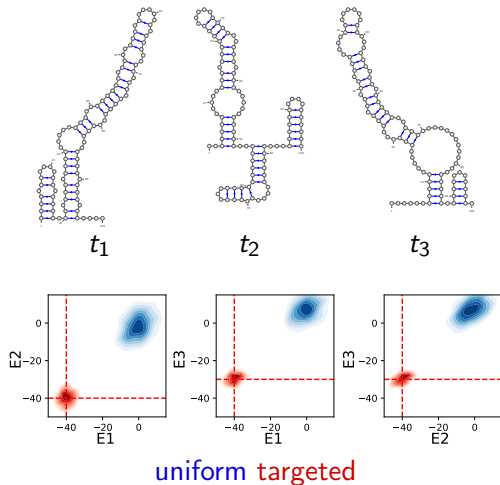
# Fixed-parameter tractable sampling in Infrared

**Recipe:**

1. **Tree-Decompose** dependency graph
2. Apply **dynamic programming** ↑ (partition functions)
3. **Sample** ↓ (stochastic traceback)

```
1 2 3 4 5 6 7

( ( . . ) ) .

. ( ( ( ) ) )

. ( ( . ) ) .
```

**target structures**

**dependency graph**

**tree decomposition**

# Fixed-parameter tractable sampling in Infrared

**Recipe:**

1. **Tree-Decompose** dependency graph
2. Apply **dynamic programming** ↑ (partition functions)
3. **Sample** ↓ (stochastic traceback)



target structures



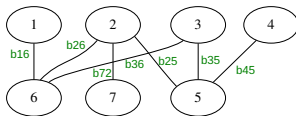dependency graph



tree decomposition

# Fixed-parameter tractable sampling in Infrared
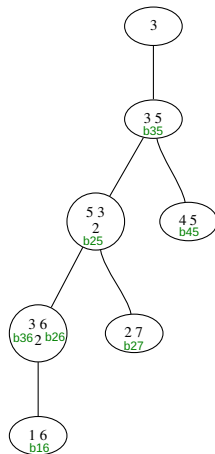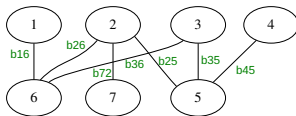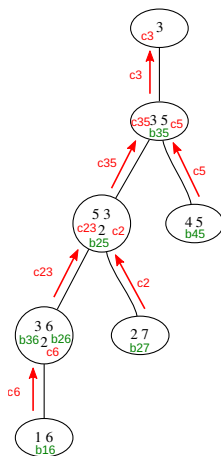
**Recipe:**

1. **Tree-Decompose** dependency graph
2. Apply **dynamic programming** $\uparrow$ (partition functions)
3. **Sample** $\downarrow$ (stochastic traceback)

```
1 2 3 4 5 6 7

( ( . . ) ) .

. ( ( ( ) ) )

. ( ( . ) ) .
```
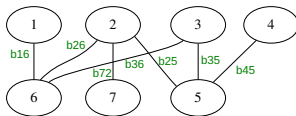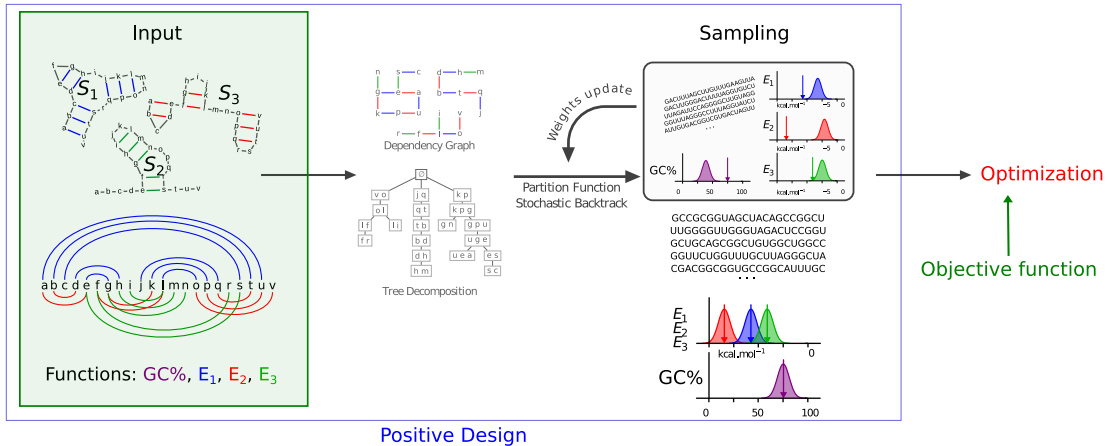
**target structures**



**dependency graph**



**tree decomposition**

**Theorem: Design sampling is efficient for fixed tree width** $w$: $\mathcal{O}(n\,k\,4^{w+1} + t\,n\,k)$

# Look back at positive design by Infrared

# RNA structure design: positive and negative

**Positive design:** Target a structure

$\rightarrow$ optimize affinity to target structures $t$

find sequence $\sigma$
with $E(\sigma, t) = \min_{\sigma'}(\sigma', t)$

extensions: multiple targets, properties, ...

**Negative design**: Avoid all off-targets

$\rightarrow$ specificity for targets

find design $\sigma$, s.t.
$E(\sigma, t) = \min_{t'} E(\sigma, t') \quad =: MFE(\sigma)$

refined objectives: probability, ensemble, ...



good positive design
since all/most other sequences
energetically worse

**-18.8**

negative design
(by MFE criterion)

**-8.2**

but no negative design, since

**-23.8**

since all other structures
energetically worse, e.g.

**-7.8**     **-6.6**     ....

# Avoiding off-targets: negative design as optimization

**minimize objective function over all sequences $\sigma$ w.r.t. a target structure $t$**

- *MFE defect:* base pair distance of MFE structure of $\sigma$ and $t$,

$$D_{MFE}(\sigma) = d(MFE(\sigma), t)$$

where *base pair distance* $d(s,t) := \sum_{(i,j)\notin s, (i,j)\in t} 1 + \sum_{(i,j)\in s, (i,j)\notin t} 1$

$$\text{Ex.: } d_{bp}\left(\begin{array}{c} (((...).).), \\ ((.(...)).) \end{array}\right) = 2$$

# Avoiding off-targets: negative design as optimization

**minimize objective function over all sequences $\sigma$ w.r.t. a target structure $t$**

- *MFE defect:* base pair distance of MFE structure of $\sigma$ and $t$,

$$D_{MFE}(\sigma) = d(MFE(\sigma), t)$$

where *base pair distance* $d(s,t) := \sum_{(i,j)\notin s, (i,j)\in t} 1 + \sum_{(i,j)\in s, (i,j)\notin t} 1$

$$\text{Ex.:} \; d_{bp} \left( \begin{array}{c} (((...).).), \\ ((.(...)).) \end{array} \right) = 2$$

- *probability defect*: $D_{Pr}(\sigma) = 1 - \Pr[t \mid \sigma]$

  maximize probability of the target $t$ in the ensemble of $\sigma$

# Avoiding off-targets: negative design as optimization

**minimize objective function over all sequences $\sigma$ w.r.t. a target structure $t$**

- *MFE defect:* base pair distance of MFE structure of $\sigma$ and $t$,

$$D_{MFE}(\sigma) = d(MFE(\sigma), t)$$

  where *base pair distance* $d(s,t) := \sum_{(i,j)\notin s,(i,j)\in t} 1 + \sum_{(i,j)\in s,(i,j)\notin t} 1$

$$\text{Ex.: } d_{bp} \left( \begin{array}{c} (((...).).), \\ ((.(...)).) \end{array} \right) = 2$$

- *probability defect*: $D_{Pr}(\sigma) = 1 - \Pr[t \mid \sigma]$
  maximize probability of the target $t$ in the ensemble of $\sigma$

  🤨 this does not consider whether the probable structures are close to target

# Avoiding off-targets: negative design as optimization

**minimize objective function over all sequences $\sigma$ w.r.t. a target structure $t$**

- *MFE defect:* base pair distance of MFE structure of $\sigma$ and $t$,
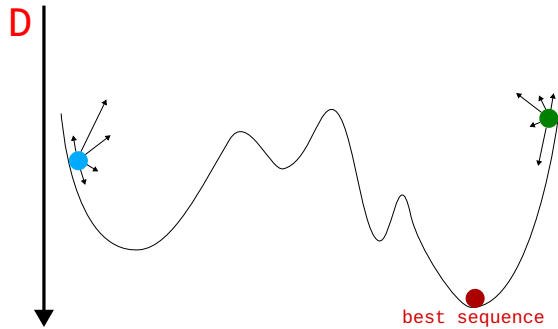
$$D_{MFE}(\sigma) = d(MFE(\sigma), t)$$

  where *base pair distance* $d(s,t) := \sum_{(i,j) \notin s, (i,j) \in t} 1 + \sum_{(i,j) \in s, (i,j) \notin t} 1$

  $$\text{Ex.: } d_{bp} \left( \begin{array}{c} (((...).).), \\ ((.(...)).) \end{array} \right) = 2$$

- *probability defect*: $D_{Pr}(\sigma) = 1 - \Pr[t \mid \sigma]$
  maximize probability of the target $t$ in the ensemble of $\sigma$

- *ensemble defect*: expected distance of ensemble structures $s$ of $\sigma$ to the target $t$

$$D_{ens}(\sigma) = \sum \Pr[s \mid \sigma] d(s,t) = \sum_{1 \le i < j \le n, (i,j) \in t} 1 - p_{ij} + \sum_{1 \le i < j \le n, (i,j) \notin t} p_{ij}$$

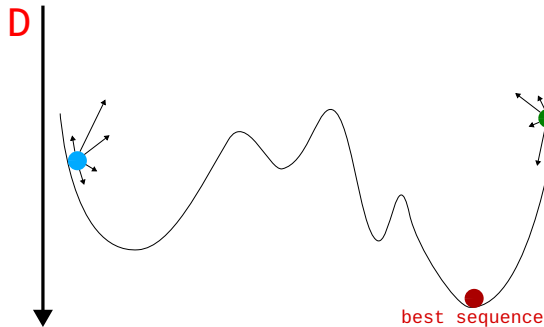# Algorithms for negative design: Stochastic Optimization



rugged landscape, local optima

random starts; neighbors, mutations

# Algorithms for negative design: Stochastic Optimization



## Hill Climbing

```
for i in range(steps):
    x = random_mutate(seq)
    if D(x) < D(seq):
        seq = x
return seq
```
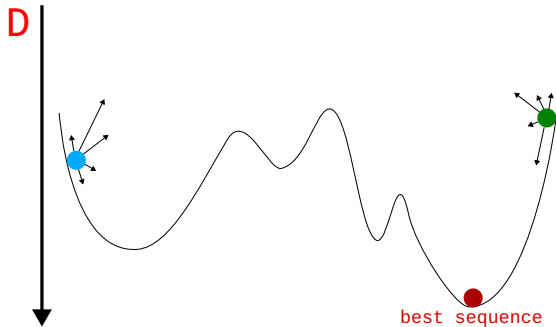
🧐 rugged landscape, local optima

🧐 random starts; neighbors, mutations

# Algorithms for negative design: Stochastic Optimization



D

best sequence

🧐 rugged landscape, local optima

🧐 random starts; neighbors, mutations

## Hill Climbing

```
for i in range(steps):
    x = random_mutate(seq)
    if D(x) < D(seq):
        seq = x
return seq
```

## Metropolis-Hastings MC Algorithm

```
best = seq
for i in range(steps):
    x = random_mutate(seq)
    if D(x) < D(seq) or
      random()<=exp((D(x)-D(seq))/T)
        seq = x
        if D(seq)<D(best): best=seq
return best
```

Control acceptance of worse neighbors by $T$
(MCMC, Boltzmann distribution)

# Algorithms for negative design: **Stochastic Optimization**



D

best sequence

🧐 rugged landscape, local optima

🧐 random starts; neighbors, mutations
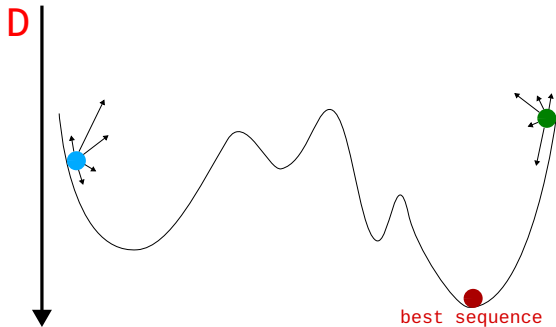
## Hill Climbing

```python
for i in range(steps):
    x = random_mutate(seq)
    if D(x) < D(seq):
        seq = x
return seq
```

## Metropolis-Hastings MC Algorithm

```python
best = seq
for i in range(steps):
    x = random_mutate(seq)
    if D(x) < D(seq) or
      random()<=exp((D(x)-D(seq))/T)
        seq = x
        if D(seq)<D(best): best=seq
return best
```

## SA, Replica exchange, Genetic algos, ...

# RNAinverse - Classical RNA design

The *single sequence design* tool of the Vienna RNA package

- Optimizes MFE defect or probability defect

[Ivo Hofacker et al., 1994]

# RNAinverse - Classical RNA design

The *single sequence design* tool of the Vienna RNA package

- Optimizes MFE defect or probability defect
- Try **random start sequences** and optimize by **Hill climbing**

needs many evaluations of objective 🧐 expensive?

[Ivo Hofacker et al., 1994]

# RNAinverse - Classical RNA design

The *single sequence design* tool of the Vienna RNA package

- Optimizes MFE defect or probability defect
- Try **random start sequences** and optimize by **Hill climbing**

  needs many evaluations of objective 🧐 expensive?

- for MFE design, RNA-tailored strategy:
  start at small substrutures, proceed to larger ones
  avoid getting stuck; reduce folding of long sequences

[Ivo Hofacker et al., 1994]

# RNAinverse - Classical RNA design

The *single sequence design* tool of the Vienna RNA package

- Optimizes MFE defect or probability defect
- Try **random start sequences** and optimize by **Hill climbing**

  needs many evaluations of objective 🤔 expensive?

- for MFE design, RNA-tailored strategy:
  start at small substrutures, proceed to larger ones
  avoid getting stuck; reduce folding of long sequences

```
( ( ( ( ( ( ( ( ( ( ( . . . ) ) ) ) ( ( ( ( . . . . ) ) ) ) ) ) ) ) ) )
? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ?
```

[Ivo Hofacker et al., 1994]

# RNAinverse - Classical RNA design

The *single sequence design* tool of the Vienna RNA package

- Optimizes MFE defect or probability defect
- Try **random start sequences** and optimize by **Hill climbing**

  needs many evaluations of objective 🧐 expensive?

- for MFE design, RNA-tailored strategy:
  start at small substrutures, proceed to larger ones
  avoid getting stuck; reduce folding of long sequences

```
( ( ( ( ( ( ( ( ( ( ( . . . ) ) ) ) ( ( ( ( . . . . ) ) ) ) ) ) ) ) )
??????CGGCAAAGCCG????????????????????
```

[Ivo Hofacker et al., 1994]

# RNAinverse - Classical RNA design

The *single sequence design* tool of the Vienna RNA package

- Optimizes MFE defect or probability defect
- Try **random start sequences** and optimize by **Hill climbing**

  needs many evaluations of objective 🧐 expensive?

- for MFE design, RNA-tailored strategy:
  start at small substrutures, proceed to larger ones
  avoid getting stuck; reduce folding of long sequences

```
( ( ( ( ( ( ( ( ( ( ( . . . ) ) ) ) ) ( ( ( ( . . . . ) ) ) ) ) ) ) ) ) )
?????CGGCAAAGCCGGGCCAAUUGGCC?????
```

[Ivo Hofacker et al., 1994]

# RNAinverse - Classical RNA design

The *single sequence design* tool of the Vienna RNA package

- Optimizes MFE defect or probability defect
- Try **random start sequences** and optimize by **Hill climbing**

  needs many evaluations of objective 🤔 expensive?

- for MFE design, RNA-tailored strategy:
  start at small substrutures, proceed to larger ones
  avoid getting stuck; reduce folding of long sequences

```
( ( ( ( ( ( ( ( ( ( ( ( . . . ) ) ) ) ) ) ( ( ( ( . . . . ) ) ) ) ) ) ) ) ) )
? ? ? ? ? CGGCAAAGCCGGGCCUUUUGGCC? ? ? ? ?
```

[Ivo Hofacker et al., 1994]

# RNAinverse - Classical RNA design

The *single sequence design* tool of the Vienna RNA package

- Optimizes MFE defect or probability defect
- Try **random start sequences** and optimize by **Hill climbing**

  needs many evaluations of objective 🧐 expensive?

- for MFE design, RNA-tailored strategy:
  start at small substrutures, proceed to larger ones
  avoid getting stuck; reduce folding of long sequences

```
( ( ( ( ( ( ( ( ( ( ( . . . ) ) ) ) ( ( ( ( . . . ) ) ) ) ) ) ) ) )
? ? ? ? ? CGGCAAAGCCGGGCCAAAUGGCC ? ? ? ? ?
```

[Ivo Hofacker et al., 1994]

# RNAinverse - Classical RNA design

The *single sequence design* tool of the Vienna RNA package

- Optimizes MFE defect or probability defect
- Try **random start sequences** and optimize by **Hill climbing**

  needs many evaluations of objective 🧐 expensive?

- for MFE design, RNA-tailored strategy:
  start at small substrutures, proceed to larger ones
  avoid getting stuck; reduce folding of long sequences
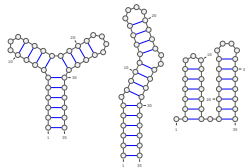
For global optima, subsequence designs must be optimal!

🧐 converse?? 🧐 why does this still work well?

RNAinverse still remarkably competitive (using good starting sequences) [Boury et al., 2024]

[Ivo Hofacker et al., 1994]

## Beyond single targets: objectives for multi-target design

```
            1         2         3
   01234567890123456789012345678901234
t1 = (((((((((((...))))(((((....))))))))))))
t2 = ((((((.(((((((((....))))...)))))))))))
t3 = .((((((...))))))).(((((((....)))))))
```



- "Multi-defect" for targets $t_1, \ldots t_m$ [Hammer et al., "RNA Blueprint" 2017]

$$
D_{multi}(\sigma) = \begin{cases} \dfrac{1}{m} \displaystyle\sum_{1 \le \ell \le m} E(\sigma, t_\ell) - G(\sigma) & \text{(dominate ensemble)} \\[3mm] + \ \xi \ \dfrac{1}{2\binom{m}{2}} \displaystyle\sum_{1 \le k < \ell \le m} |E(\sigma, t_k) - E(\sigma, t_\ell)| & \text{(similar energies)} \end{cases}
$$

$G(\sigma) = -RT \ln Z(\sigma)$ "ensemble energy"

## Beyond single targets: objectives for multi-target design

```
          1         2         3
 01234567890123456789012345678901234
t1 = (((((((((((...))))(((((....)))))))))))
t2 = (((((((.(((((((((....))))))..)))))))))
t3 = .((((((...))))))).(((((((....)))))))
```



- "Multi-defect" for targets $t_1, \ldots t_m$ [Hammer et al., "RNA Blueprint" 2017]

$$
D_{multi}(\sigma) = \left\{
\begin{array}{ll}
\dfrac{1}{m} \displaystyle\sum_{1 \leq \ell \leq m} E(\sigma, t_\ell) - G(\sigma) & \text{(dominate ensemble)} \\[2mm]
+ \;\; \xi \; \dfrac{1}{2\binom{m}{2}} \displaystyle\sum_{1 \leq k < \ell \leq m} |E(\sigma, t_k) - E(\sigma, t_\ell)| & \text{(similar energies)}
\end{array}
\right.
$$

- Aim for ensemble dominance with certain energy differences of targets..., e.g.

$$
D_{ex}(\sigma) = |E(\sigma, t_1) - G(\sigma)| + |E(\sigma, t_2) - E(\sigma, t_1) - 3| + |E(\sigma, t_3) - E(\sigma, t_1) - 4|
$$

$$
G(\sigma) = -RT \ln Z(\sigma) \quad \text{"ensemble energy"}
$$

# Stochastic optimization in Infrared

🧐 How to find valid neighbors in complex design problems?

```
          1                 2               3
0123456789012345678901234567890123456789012345678901234
( ( ( ( ( ( ( ( ( ( . . . ) ) ) ) ( ( ( ( . . . . ) ) ) ) ) ) ) ) ) )
( ( ( ( ( ( . ( ( ( ( ( ( ( ( . . . . ) ) ) ) . . ) ) ) ) ) ) ) ) ) )
. ( ( ( ( ( ( . . . ) ) ) ) ) ) . ( ( ( ( ( ( ( . . . . ) ) ) ) ) ) )
GCGUGCGGGGGAGUCUCUCCGUCAAUGGGGCACGC
```

# Stochastic optimization in Infrared

🧐 How to find valid neighbors in complex design problems?

```
          1                   2                   3
012345678901234567890123456789012345678901234
((((((((((...))))((((....))))))))))
((((((.((((((((....))))..))))))))))
.((((((...))))))).(((((((....)))))))
GCGUGCGGGGGAGUCUCUCCGUCAAUGGGGCACGC
```



- resample connected components (of dependency graph)
    Idea: independence of cc preserves all other constraints

# Stochastic optimization in Infrared

🧐 How to find valid neighbors in complex design problems?

```
          1         2         3
012345678901234567890123456789012 34
((((((((((...))))(((( ....))))))))))
(((((( .((((((((( ....)))) .. )))))))))))
 .((((((( ...)))))) .((((((( .... )))))))
GCGUGCGGGGGAGUCUCUCCGUCAAUGGGGCACGC
```



- resample connected components (of dependency graph)

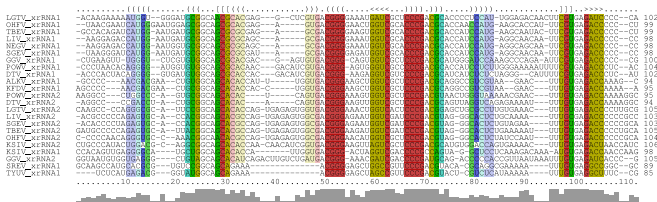    Idea: independence of cc preserves all other constraints

- sample neighbors in a targeted distance

    Idea: • extend model by distance function
         • sample neighbors, controlled by distance to current sequence

# Learning design from evolution (Generative Models)

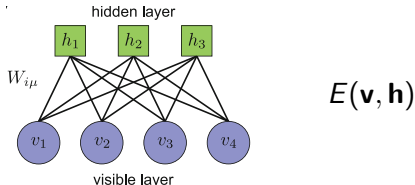General idea: learn from homology information / MSAs



*There is information in MSA beyond position-wise frequencies!* $\Rightarrow$ covariation...

$\Rightarrow$ Restricted Boltzmann Machines (RBM)

- (bipartite) two layer neural networks
- can be trained to evaluate sequences
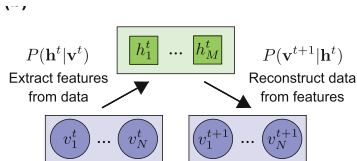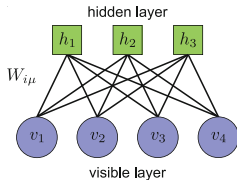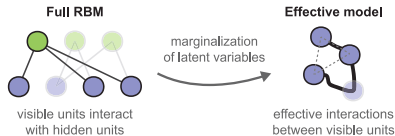- shown to design SAM aptamer
  [FdCD et al., 2023]



$E(\mathbf{v}, \mathbf{h})$

[Jorge Fernandez-de-Cossio-Diaz, 2024]

RNA Design · S. Will · 29

# Restricted Boltzmann Machines (RBMs)



$v_\mu$: A, C, G, U, –
$h_i$: dependencies

hidden layer

$W_{i\mu}$

visible layer

$P(\mathbf{h}^t | \mathbf{v}^t)$
Extract features from data

$P(\mathbf{v}^{t+1} | \mathbf{h}^t)$
Reconstruct data from features

PCD, Gibbs sampler

**Full RBM**

**Effective model**

marginalization of latent variables

visible units interact with hidden units

effective interactions between visible units
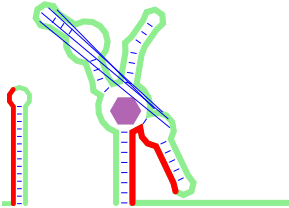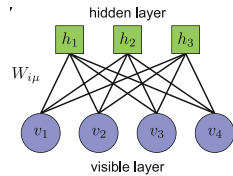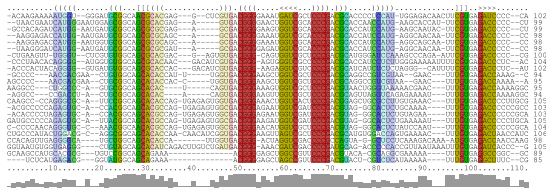
$E_{eff}(\mathbf{v}) \leftarrow$ marginalization

$$E(\mathbf{v}, \mathbf{h}) = \sum_{i=1}^{N} \mathcal{V}_i(v_i) + \sum_{\mu=1}^{M} \mathcal{U}_\mu(h_\mu) - \sum_{i=1}^{N} \sum_{\mu=1}^{M} w_{i\mu}(v_i) h_\mu$$

- Effective training: maximize log likelihood of training data by *persistent contrastive divergence (PDC)* [Hinton, 2012] • Evaluation of designs • Positive design / sampling

[Jorge Fernandez-de-Cossio-Diaz, 2024]

# Hands on session

- Use of Infrared
- Design of SAM-I aptamer
- Design of SAM-I on-switch (simplified)

- Sample and optimize
- Integrate homology and thermodynamic info
- Integrate evaluation by RBM