# Metropolis Monte Carlo sampling

# Gibbs Clustering

## Morten Nielsen

## Department of Health Technology
## DTU, Denmark

# Metropolis Monte Carlo

- What to do if you cannot do the math or your error function cannot be differentiated?

$$\frac{\partial E}{\partial w_i} = \ ?$$

# Example: Estimating Π by Independent Monte-Carlo Samples

Suppose we throw darts randomly (and uniformly) at the square:

$$\frac{\#\ \text{darts hitting shaded area}}{\#\ \text{darts hitting inside square}} = \frac{\frac{1}{4}\pi r^2}{r^2} = \frac{1}{4}\pi$$

or

$$\pi = 4\ \frac{\#\ \text{darts hitting shaded area}}{\#\ \text{darts hitting inside square}}$$

**Algorithm:**
**For i=[1..ntrials]**
    **x = (random# in [0..r])**
    **y = (random# in [0..r])**
    **distance = sqrt (x^2 + y^2)**
    **if distance ≤ r**
       **hits++**
**End**
**Output:** $\dfrac{4 \times hits}{ntrials}$



http://www.chem.unl.edu/zeng/joy/mclab/mcintro.html

# Estimating Π

# Monte Carlo

Because of their reliance on repeated computation of random or pseudo-random numbers, Monte Carlo methods are most suited to calculation by a computer. Monte Carlo methods tend to be used when it is unfeasible or impossible to compute an exact result with a deterministic algorithm
Or when you are too stupid to do the math yourself?

$$E = f(x)$$

$$dE = E_1 - E_0$$

$$P(accept) = \min(1, e^{-dE/T})$$

# Class II MHC binding

- MHC class II binds peptides in the class II antigen presentation pathway
- Binds peptides of length 9-18 (even whole proteins can bind!)
- Binding cleft is open
- Binding core is 9 aa

# Conventional Gibbs sampling
# MHC class II binding

9 AAs

SLFIGLKGDIRESTV
DGEEI
VFRLI
SFSC
IDQV
WIQKI
KMLLI
ELLEI
LNKF
ESLHI

DFAA

SLFIGLKGDIRESTV
DGEEI
VFRL

IDQ
WIQKE
KMLL

ESLHN

DFA

$$E = \sum_{peptides} \log \frac{p_{p,a}}{q_a}$$

# Gibbs sampling - sequence alignment

CENTERFO
RBIOLOGI
CALSEQU
ENCEANA
LYSIS CBS

CENTER FOR BIOLOGICAL SEQUENCE ANALYSIS

## Why sampling?

50 sequences 12 amino acids long

try all possible combinations with a 9-mer overlap

$4^{50} \sim 10^{30}$ possible combinations

...computationally unfeasible

```
 SLFIGLKGDIRESTV
DGEEVQLIAAVPGK
 VFRLKGGAPIKGVTF
     SFSCIAIGIITLYLG
  IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
 KMLLDNINTPEGIIP
      ELLEFHYYLSSKLNK
     LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
        . . .
        . . .
        . . .
        . . .
  DFAAQVDYPSTGLY
```

# Gibbs sampling - sequence alignment

## State transition

```
   SLFIGLKGDIRESTV
DGEEEVQLIAAVPGK
  VFRLKGGAPIKGVTF
     SFSCIAIGIITLYLG
   IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
  KMLLDNINTPEGIIP
      ELLEFHYYLSSKLNK
     LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
```

move to
state +1

$$E = \sum_{peptides} \log \frac{p_{p,a}}{q_a}$$

$$dE = E_i - E_{i-1}$$

# Gibbs sampling - sequence alignment

## State transition



move to
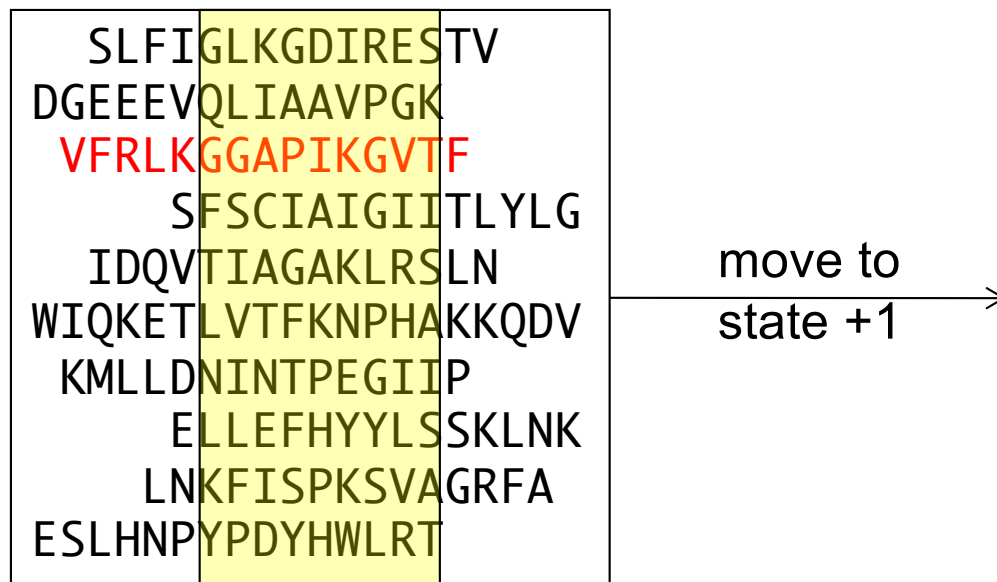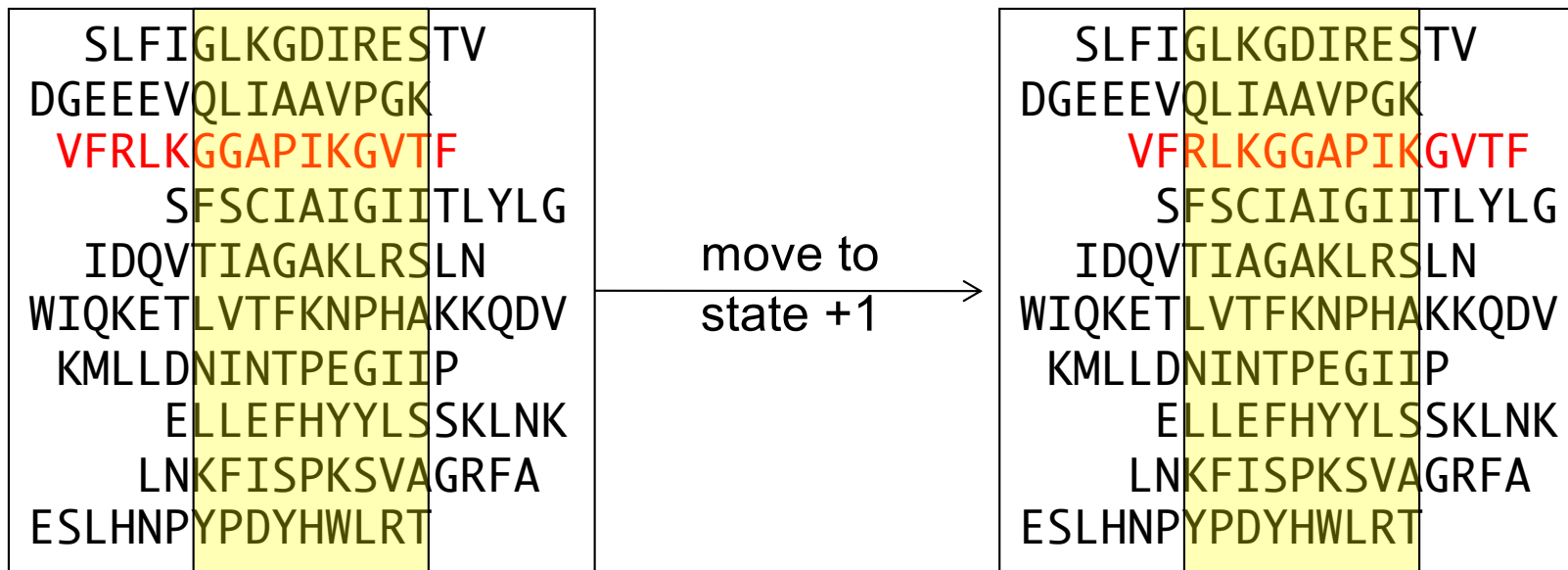state +1

```
  SLFIGLKGDIRESTV
DGEEEVQLIAAVPGK
 VFRLKGGAPIKGVTF
     SFSCIAIGIITLYLG
   IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
  KMLLDNINTPEGIIP
       ELLEFHYYLSSKLNK
     LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
```

```
  SLFIGLKGDIRESTV
DGEEEVQLIAAVPGK
    VFRLKGGAPIKGVTF
     SFSCIAIGIITLYLG
   IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
  KMLLDNINTPEGIIP
       ELLEFHYYLSSKLNK
     LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
```

$$E = \sum_{peptides} \log \frac{p_{p,a}}{q_a}$$

$$dE = E_i - E_{i-1}$$

Accept or reject the move?

$$P = \min\left[1, \exp\left(\frac{dE}{T}\right)\right]$$

**DTU**

Department of Systems Biology
Technical University of Denmark

Note that the probability of going to the new
state depends on the previous state only

# Gibbs sampling - sequence alignment

## Numerical example - 1



```
  SLFIGLKGDIRESTV
DGEEEVQLIAAVPGK
 VFRLKGGAPIKGVTF
     SFSCIAIGIITLYLG
  IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
  KMLLDNINTPEGIIP
      ELLEFHYYLSSKLNK
    LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
```

move to
state +1

$T = 0.2$

```
  SLFIGLKGDIRESTV
DGEEEVQLIAAVPGK
 VFRLKGGAPIKGVTF
     SFSCIAIGIITLYLG
  IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
  KMLLDNINTPEGIIP
      ELLEFHYYLSSKLNK
    LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
```

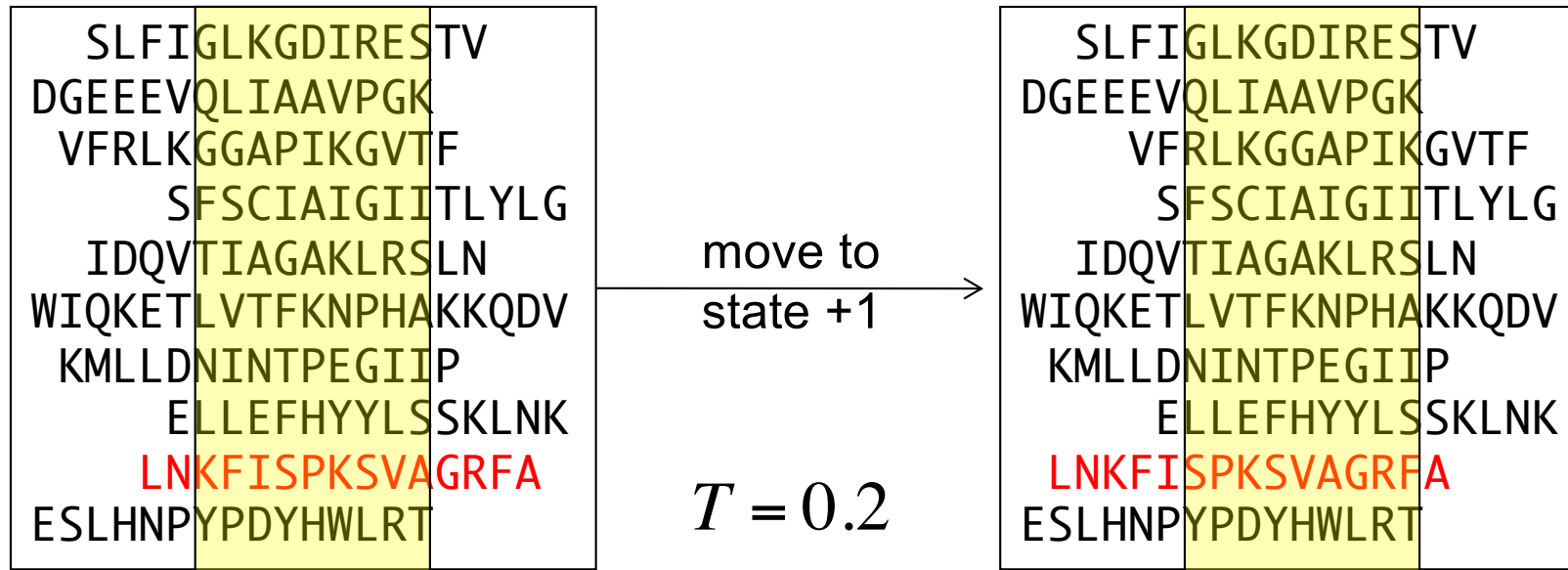$$E_{i-1} = 2.44 \qquad\qquad E_i = 2.52$$

$$P = \min\left[1, \exp\left(\frac{0.08}{0.2}\right)\right] = \min[1 \,,\, 1.49] = 1$$

**Accept move with Prob = 100%**

# Gibbs sampling - sequence alignment

## Numerical example - 2

```
   SLFIGLKGDIRESTV
DGEEEVQLIAAVPGK
  VFRLKGGAPIKGVTF
      SFSCIAIGIITLYLG
  IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
  KMLLDNINTPEGIIP
      ELLEFHYYLSSKLNK
    LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
```

move to
state +1

$T = 0.2$

```
   SLFIGLKGDIRESTV
DGEEEVQLIAAVPGK
  VFRLKGGAPIKGVTF
      SFSCIAIGIITLYLG
  IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
  KMLLDNINTPEGIIP
      ELLEFHYYLSSKLNK
    LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
```
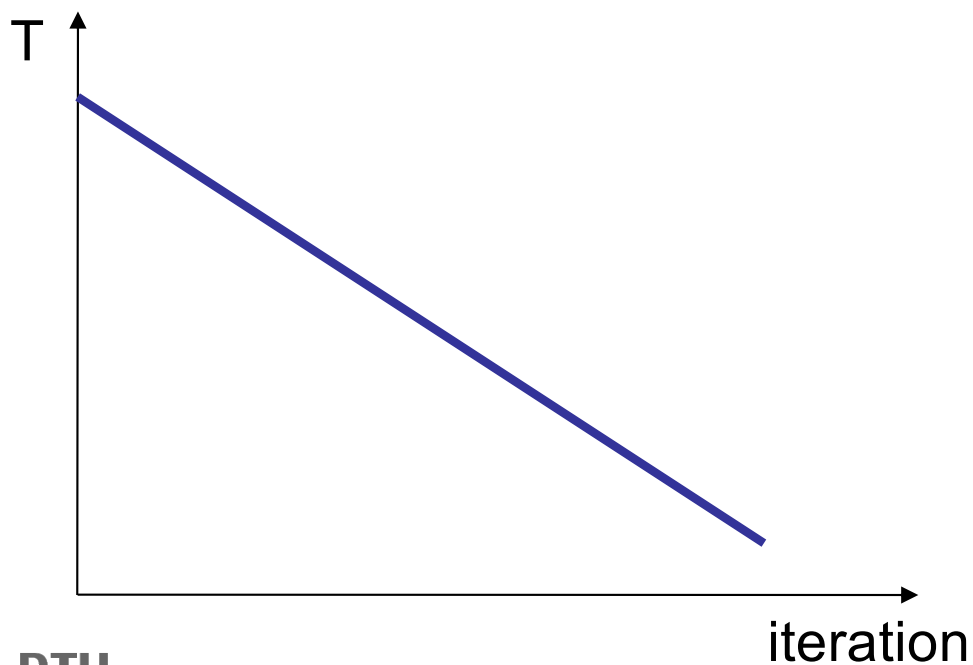
$$E_{i-1} = 2.44$$

$$E_i = 2.35$$

$$P = \min\left[1, \exp\left(\frac{-0.09}{0.2}\right)\right] = \min\left[1 , 0.638\right] = 0.638$$

**Accept move with Prob = 63.8%**

# Gibbs sampling - sequence alignment

## What is the MC temperature?

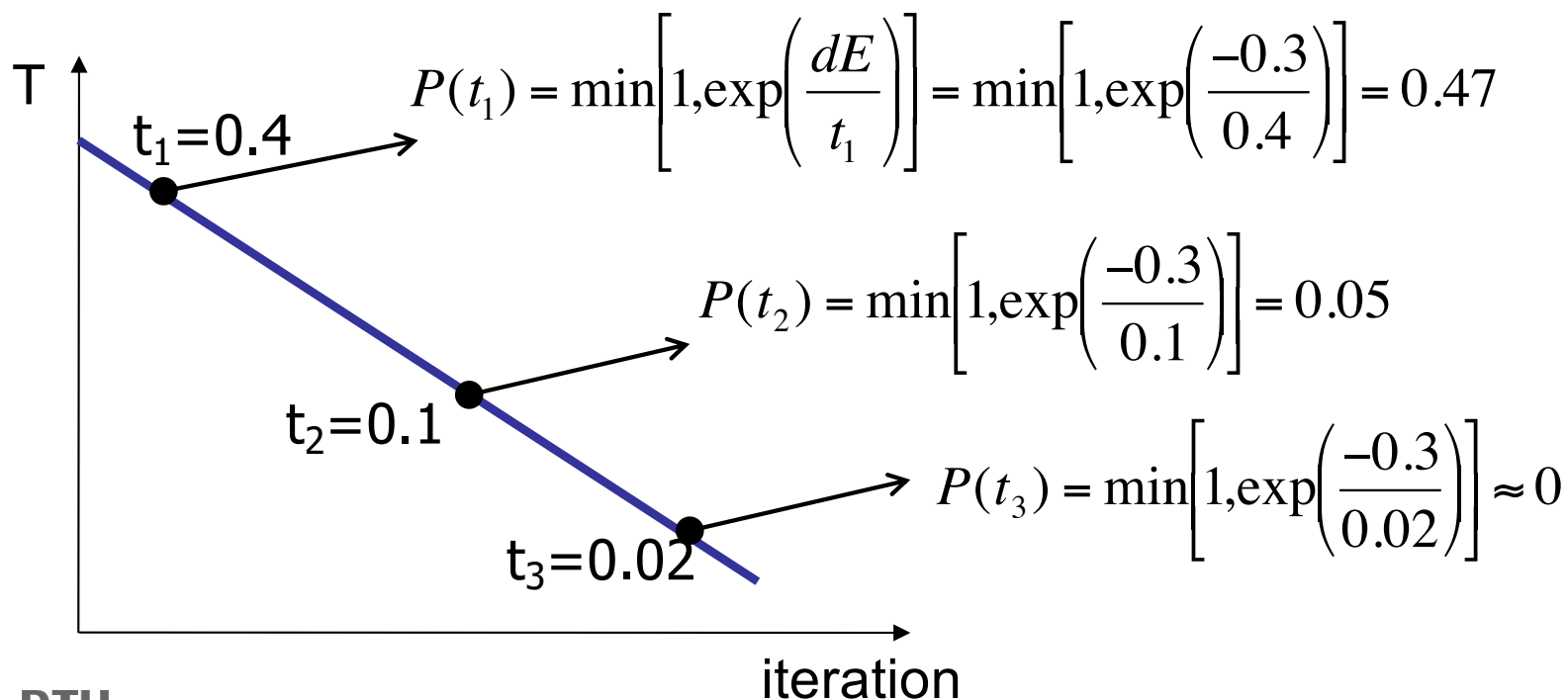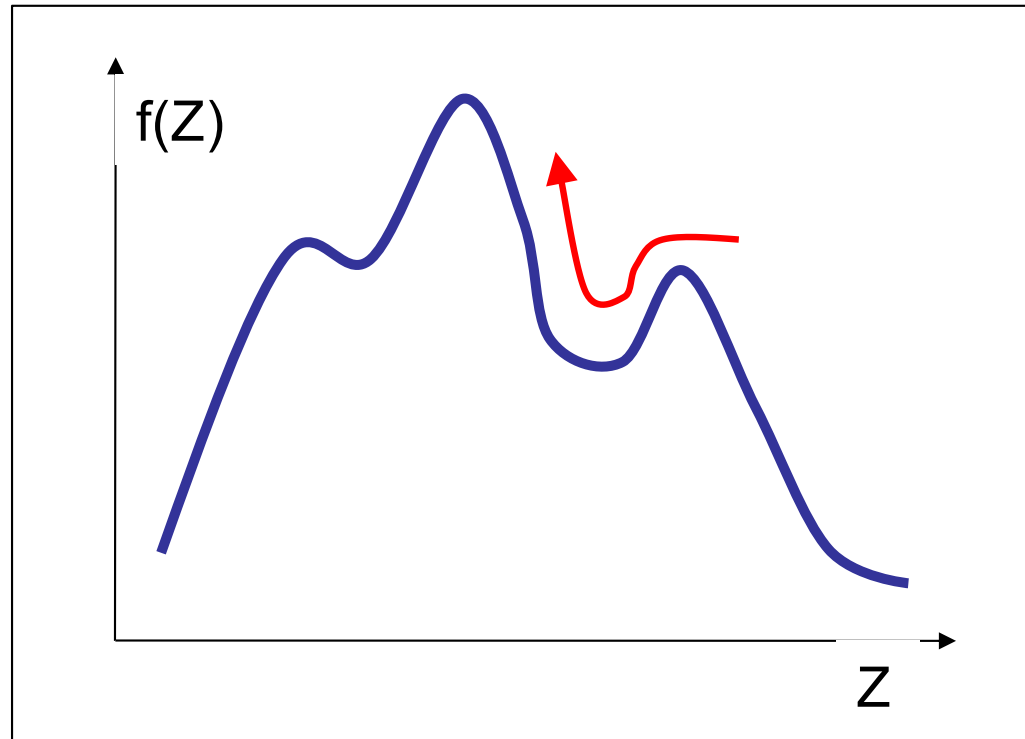it's a scalar decreased during the simulation

# Gibbs sampling - sequence alignment

## What is the MC temperature?

it's a scalar decreased during the simulation

E.g. same **dE=-0.3** but at different temperatures

$$P(t_1) = \min\left[1, \exp\left(\frac{dE}{t_1}\right)\right] = \min\left[1, \exp\left(\frac{-0.3}{0.4}\right)\right] = 0.47$$

$$P(t_2) = \min\left[1, \exp\left(\frac{-0.3}{0.1}\right)\right] = 0.05$$

$$P(t_3) = \min\left[1, \exp\left(\frac{-0.3}{0.02}\right)\right] \approx 0$$

T

$t_1 = 0.4$

$t_2 = 0.1$

$t_3 = 0.02$

iteration

Move freely around states when the system is "warm", then cool it off to force it into a state of high fitness

# Local minima

DTU

Department of Syst
Technical University of Denmark

# Does it work?

CENTERFO
RBIOLOGI
CALSEQU
ENCEANA
LYSIS CBS

CENTER FOR BIOLOGICAL SEQUENCE ANALYSIS

HLA-DRB3*01:01



created by Seq2Logo

DTU

Department of Systems Biology
Technical University of Denmark

# It works

CENTERFO
RBIOLOGI
CALSEQU
ENCEANA
LYSIS **CBS**

CENTER FOR BIOLOGICAL SEQUENCE ANALYSIS

# Gibbs sampler. Prediction accuracy



| | Gibbs | Tepitope | ClustalW |
|---|---|---|---|
| AUC | 0.744 | 0.702 | 0.667 |

# Interpreting and benefitting from MS eluted ligand data sets

A: Sample preparation

Cells lysate — Immuno-affinity purification — HLA-I complexes — Purification and enrichment — HLA-I peptides

*Michel Bassani-Sternberg et al, MCP, 2015*

*Gibbs Cluster analysis*



Input Peptides

Clustered Output

Group 3, Group 2, Group 1

Information Content

# Groups

*GibbsCluster, Andreatta, Alvarez, Nielsen, NAR, 2017*

# The algorithm

CENTERFO
RBIOLOGI
CALSEQU
ENCEANA
LYSIS CBS

CENTER FOR BIOLOGICAL SEQUENCE ANALYSIS

### 1. List of peptides

```
        SLFIGLKGDIRESTV
        DGEEEVQLIAAVPGK
        VFRLKGGAPIKGVTF
        SFSCIAIGIITLYLG
        IDQVTIAGAKLRSLN
    WIQKETLVTFKNPHAKKQDV
        KMLLDNINTPEGIIP
        ELLEFHYYLSSKLNK
        LNKFISPKSVAGRFA
        ESLHNPYPDYHWLRT
        NKVKSLRILNTRRKL
        MMGMFNMLSTVLGVS
        AKSSPAYPSVLGQTI
      RHLIFCHSKKKCDELAAK
```

# The algorithm

**1. List of peptides**

```
SLFIGLKGDIRESTV
DGEEEVQLIAAVPGK
VFRLKGGAPIKGVTF
SFSCIAIGIITLYLG
IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
KMLLDNINTPEGIIP
ELLEFHYYLSSKLNK
LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
NKVKSLRILNTRRKL
MMGMFNMLSTVLGVS
AKSSPAYPSVLGQTI
RHLIFCHSKKKCDELAAK
```

**2. create N random groups**

**g₁**

```
----IDQVTIAGAKLRSLN-
WIQKETLVTFKNPHAKKQDV
--ELLEFHYYLSSKLNK---
--MMGMFNMLSTVLGVS---
---AKSSPAYPSVLGQTI--
```

**g₂**

```
-SLFIGLKGDIRESTV--
---SFSCIAIGIITLYLG
KMLLDNINTPEGIIP---
-LNKYVHGTWRSILP---
--NKVKSLRILNTRRKL-
```

**gₙ**

```
---ESLHNPYPDYHWLRT-
RHLIFCHSKKKCDELAAK-
----VFRLKGGAPIKGVTF
--LNKFISPKSVAGRFA--
DGEEEVQLIAAVPGK----
```

**3     Simple shift     Remove peptide**

# The algorithm

CENTERFOR
BIOLOGI
CALSEQU
ENCEANA
LYSIS **CBS**

CENTER FOR BIOLOGICAL SEQUENCE ANALYSIS

**1. List of peptides**

```
SLFIGLKGDIRESTV
DGEEEVQLIAAVPGK
VFRLKGGAPIKGVTF
SFSCIAIGIITLYLG
IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
KMLLDNINTPEGIIP
ELLEFHYYLSSKLNK
LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
NKVKSLRILNTRRKL
MMGMFNMLSTVLGVS
AKSSPAYPSVLGQTI
RHLIFCHSKKKCDELAAK
```
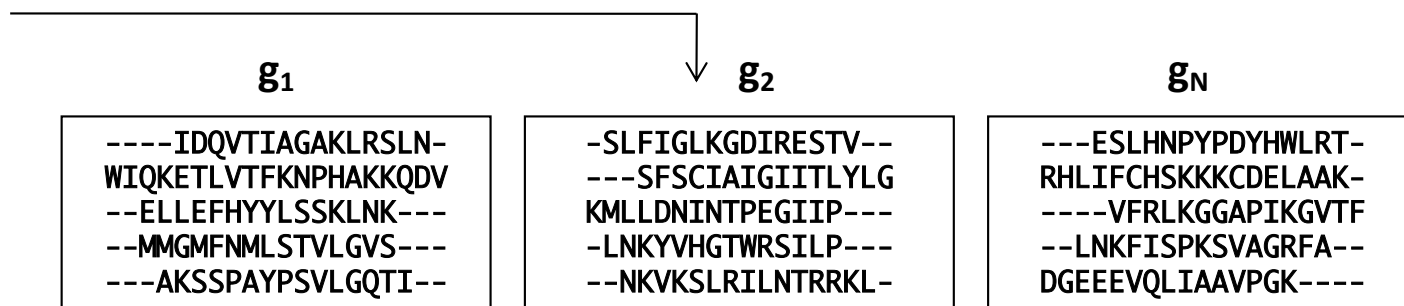
**2. create N random groups**

**g₁**

```
----IDQVTIAGAKLRSLN-
WIQKETLVTFKNPHAKKQDV
--ELLEFHYYLSSKLNK---
--MMGMFNMLSTVLGVS---
---AKSSPAYPSVLGQTI--
```

**g₂**

```
-SLFIGLKGDIRESTV--
---SFSCIAIGIITLYLG
KMLLDNINTPEGIIP---
-LNKYVHGTWRSILP---
--NKVKSLRILNTRRKL-
```

**gₙ**

```
---ESLHNPYPDYHWLRT-
RHLIFCHSKKKCDELAAK-
----VFRLKGGAPIKGVTF
--LNKFISPKSVAGRFA--
DGEEEVQLIAAVPGK----
```

**3. Finding the optimal configuration (the MC moves)**

**MMGMFNMLSTVLGVS**

# The algorithm

CENTERFO
RBIOLOGI
CALSEQU
ENCEANA
LYSIS **CBS**

CENTER FOR BIOLOGICAL SEQUENCE ANALYSIS

**1. List of peptides**

```
SLFIGLKGDIRESTV
DGEEEVQLIAAVPGK
VFRLKGGAPIKGVTF
SFSCIAIGIITLYLG
IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
KMLLDNINTPEGIIP
ELLEFHYYLSSKLNK
LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
NKVKSLRILNTRRKL
MMGMFNMLSTVLGVS
AKSSPAYPSVLGQTI
RHLIFCHSKKKCDELAAK
```
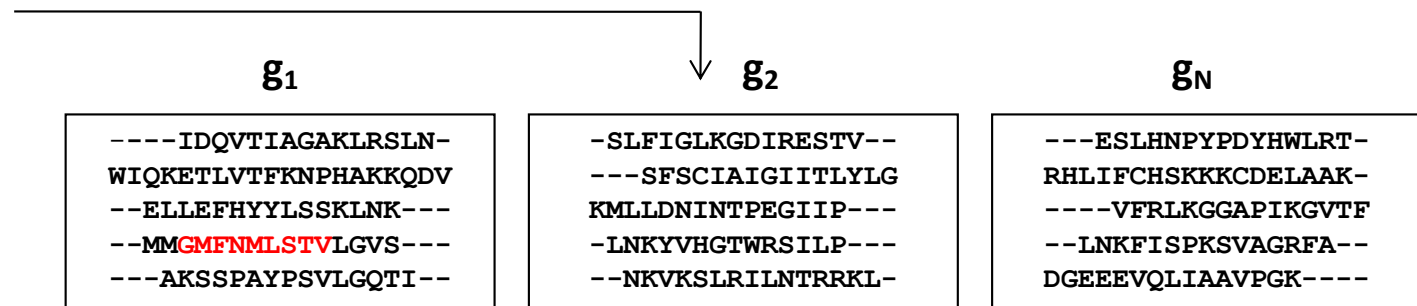
**2. create N
random groups**

**g₁**

```
----IDQVTIAGAKLRSLN-
WIQKETLVTFKNPHAKKQDV
--ELLEFHYYLSSKLNK---
---AKSSPAYPSVLGQTI--
```

**g₂**

```
-SLFIGLKGDIRESTV--
---SFSCIAIGIITLYLG
KMLLDNINTPEGIIP---
-LNKYVHGTWRSILP---
--NKVKSLRILNTRRKL-
```

**gₙ**

```
---ESLHNPYPDYHWLRT-
RHLIFCHSKKKCDELAAK-
----VFRLKGGAPIKGVTF
--LNKFISPKSVAGRFA--
DGEEEVQLIAAVPGK----
```

**3. Finding the optimal configuration (the MC moves)**

MMGMFNMLSTVLGVS

**5. Score new core to log-
odds matrices**

**4. Random shift of the core**

MMGMFNMLSTVLGVS

$$dE = E_{new} - E_{old}$$

*6. Accept or reject
move*

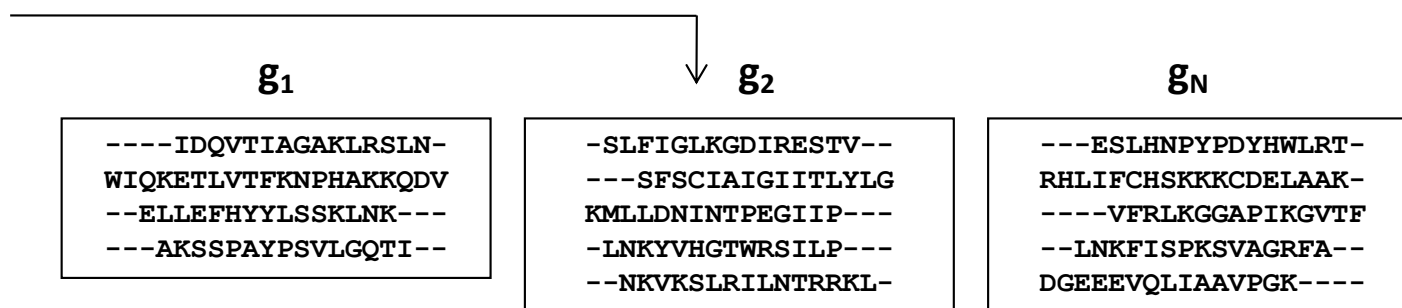$$P = \min\left[1, \exp\left(\frac{dE}{T}\right)\right]$$

# *The* algorithm

**1. List of peptides**

```
SLFIGLKGDIRESTV
DGEEEVQLIAAVPGK
VFRLKGGAPIKGVTF
SFSCIAIGIITLYLG
IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
KMLLDNINTPEGIIP
ELLEFHYYLSSKLNK
LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
NKVKSLRILNTRRKL
MMGMFNMLSTVLGVS
AKSSPAYPSVLGQTI
RHLIFCHSKKKCDELAAK
```

**2. create N random groups**

**g₁**

```
----IDQVTIAGAKLRSLN-
WIQKETLVTFKNPHAKKQDV
--ELLEFHYYLSSKLNK---
--MMGMFNMLSTVLGVS---
---AKSSPAYPSVLGQTI--
```

**g₂**

```
-SLFIGLKGDIRESTV--
---SFSCIAIGIITLYLG
KMLLDNINTPEGIIP---
-LNKYVHGTWRSILP---
--NKVKSLRILNTRRKL-
```

**gₙ**

```
---ESLHNPYPDYHWLRT-
RHLIFCHSKKKCDELAAK-
----VFRLKGGAPIKGVTF
--LNKFISPKSVAGRFA--
DGEEEVQLIAAVPGK----
```

**3** **Simple**        **Remove peptide**

**MMGMFNMLSTVLGVS**

# The algorithm

**1. List of peptides**

```
SLFIGLKGDIRESTV
DGEEEVQLIAAVPGK
VFRLKGGAPIKGVTF
SFSCIAIGIITLYLG
IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
KMLLDNINTPEGIIP
ELLEFHYYLSSKLNK
LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
NKVKSLRILNTRRKL
MMGMFNMLSTVLGVS
AKSSPAYPSVLGQTI
RHLIFCHSKKKCDELAAK
```

**2. create N random groups**

**$g_1$**

```
----IDQVTIAGAKLRSLN-
WIQKETLVTFKNPHAKKQDV
--ELLEFHYYLSSKLNK---
---AKSSPAYPSVLGQTI--
```

**$g_2$**

```
-SLFIGLKGDIRESTV--
---SFSCIAIGIITLYLG
KMLLDNINTPEGIIP---
-LNKYVHGTWRSILP---
--NKVKSLRILNTRRKL-
```
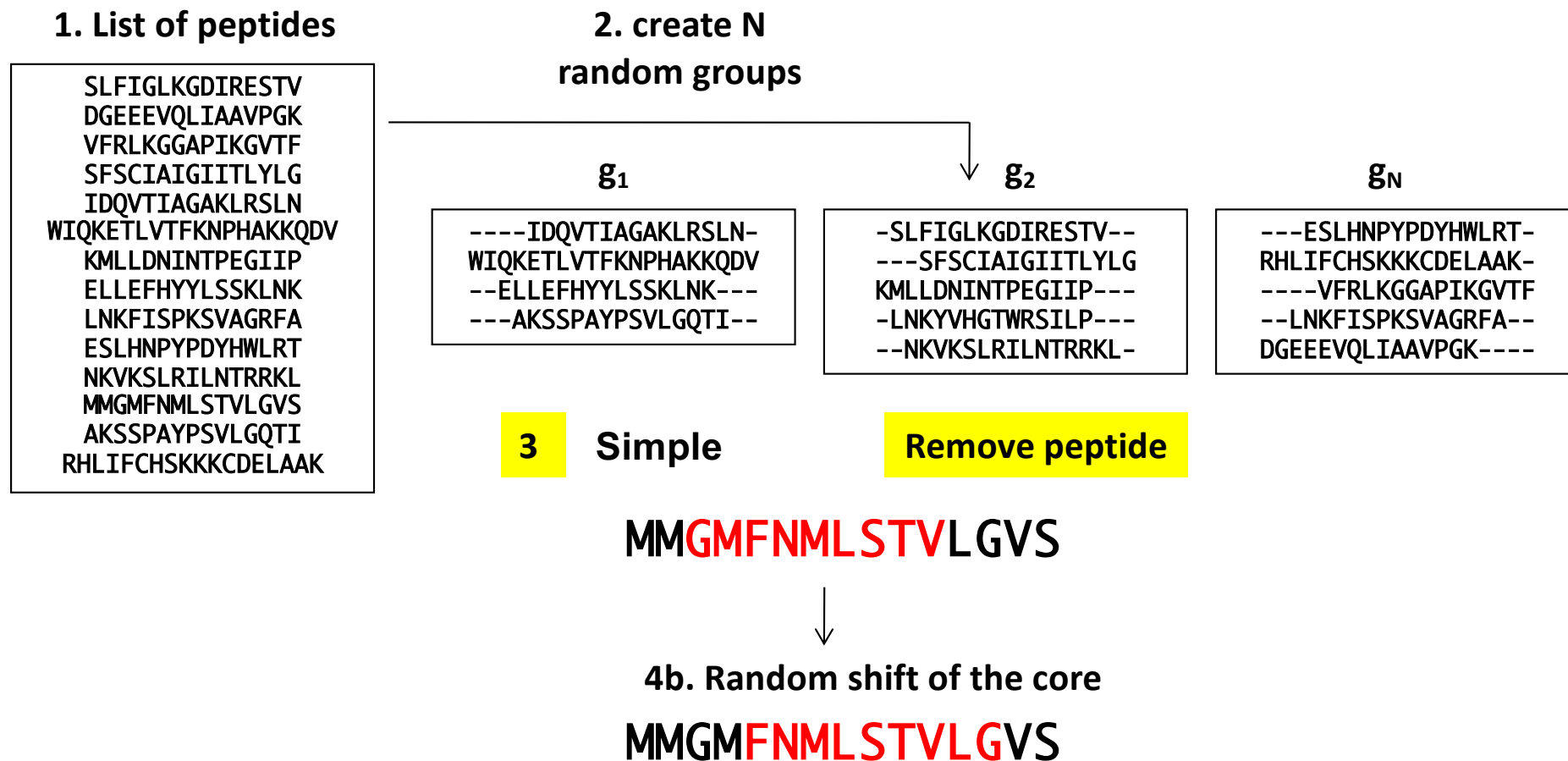
**$g_N$**

```
---ESLHNPYPDYHWLRT-
RHLIFCHSKKKCDELAAK-
----VFRLKGGAPIKGVTF
--LNKFISPKSVAGRFA--
DGEEEVQLIAAVPGK----
```
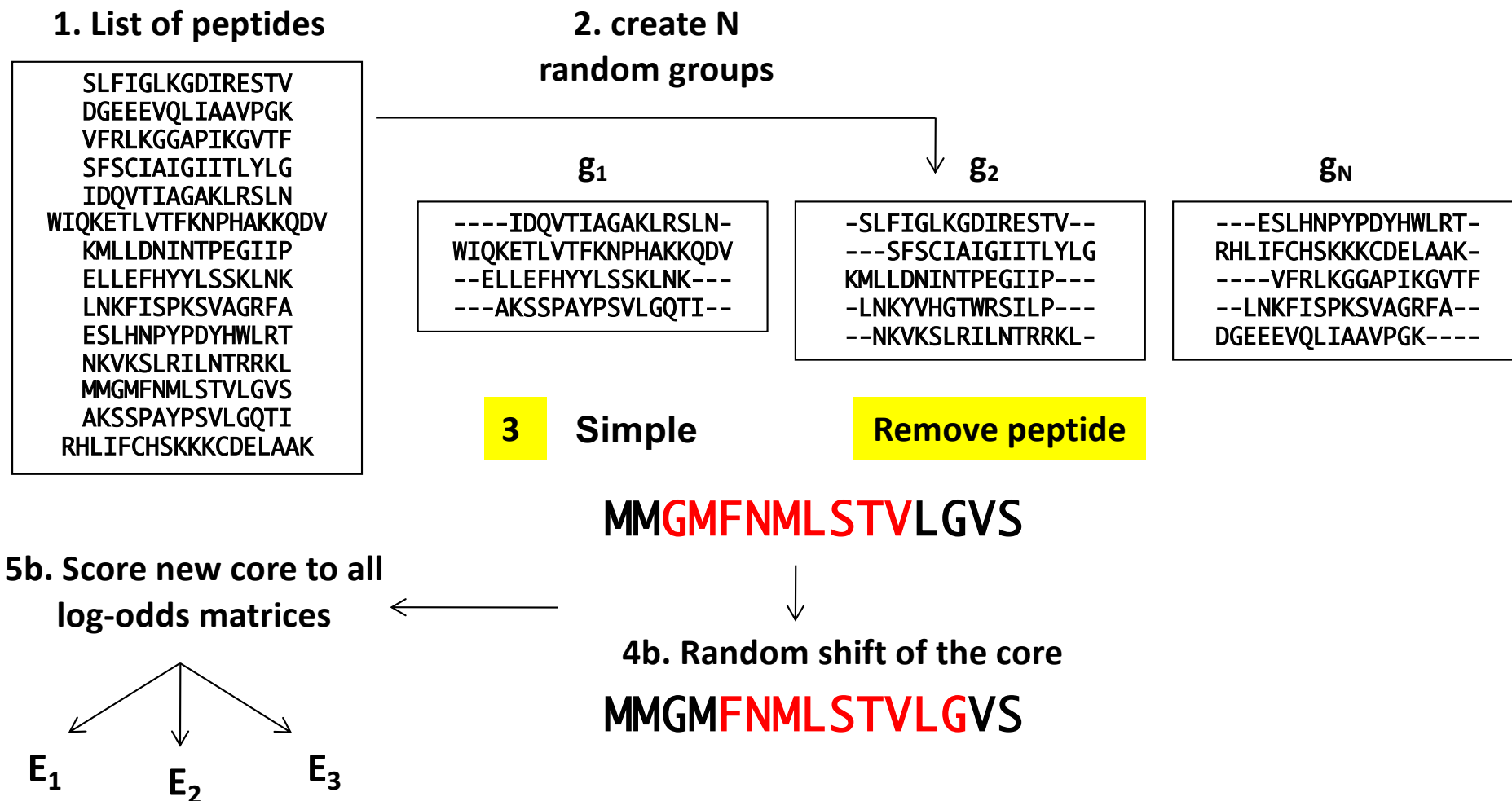
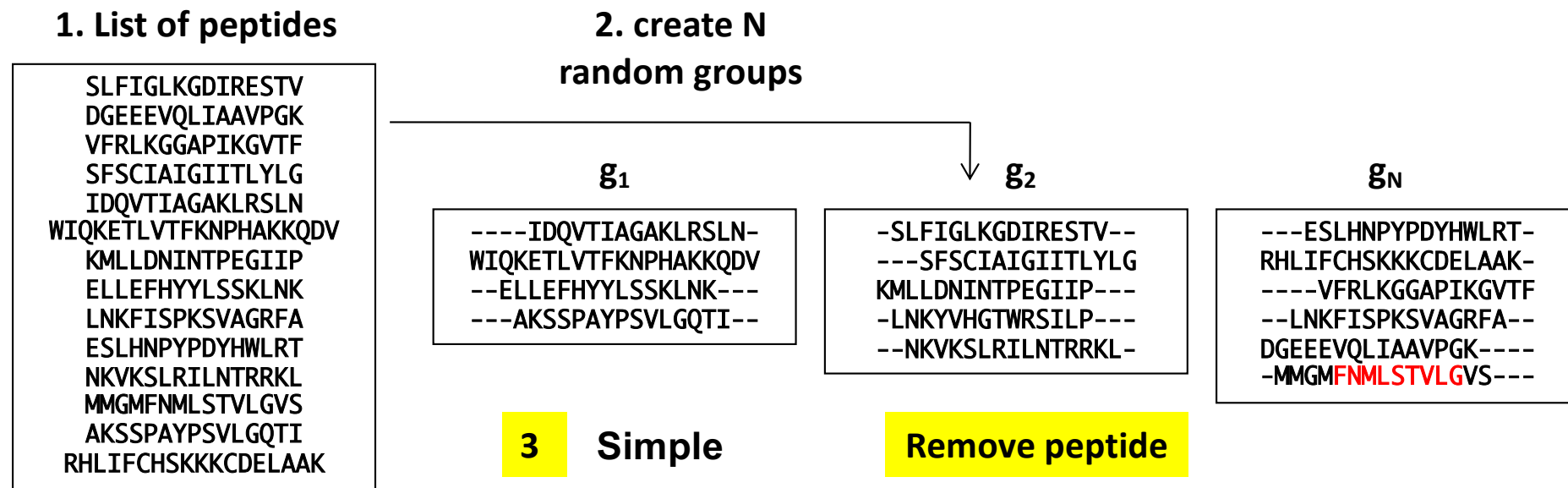**3** **Simple** **Remove peptide**

MMGMFNMLSTVLGVS

↓

**4b. Random shift of the core**

MMGMFNMLSTVLGVS

# The algorithm

**1. List of peptides**

```
SLFIGLKGDIRESTV
DGEEEVQLIAAVPGK
VFRLKGGAPIKGVTF
SFSCIAIGIITLYLG
IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
KMLLDNINTPEGIIP
ELLEFHYYLSSKLNK
LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
NKVKSLRILNTRRKL
MMGMFNMLSTVLGVS
AKSSPAYPSVLGQTI
RHLIFCHSKKKCDELAAK
```

**2. create N random groups**

**g₁**

```
----IDQVTIAGAKLRSLN-
WIQKETLVTFKNPHAKKQDV
--ELLEFHYYLSSKLNK---
---AKSSPAYPSVLGQTI--
```

**g₂**

```
-SLFIGLKGDIRESTV--
---SFSCIAIGIITLYLG
KMLLDNINTPEGIIP---
-LNKYVHGTWRSILP---
--NKVKSLRILNTRRKL-
```

**g_N**

```
---ESLHNPYPDYHWLRT-
RHLIFCHSKKKCDELAAK-
----VFRLKGGAPIKGVTF
--LNKFISPKSVAGRFA--
DGEEEVQLIAAVPGK----
```

**3** **Simple** **Remove peptide**

MMGMFNMLSTVLGVS

**5b. Score new core to all log-odds matrices**

**4b. Random shift of the core**

MMGMFNMLSTVLGVS

E₁  E₂  E₃

$$dE = E_{before} - \max(E_1, E_2, E_3)$$

# The algorithm

CENTERFO
RBIOLOGI
CALSEQU
ENCEANA
LYSIS **CBS**

CENTER FOR BIOLOGICAL SEQUENCE ANALYSIS

**1. List of peptides**

```
SLFIGLKGDIRESTV
DGEEEVQLIAAVPGK
VFRLKGGAPIKGVTF
SFSCIAIGIITLYLG
IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
KMLLDNINTPEGIIP
ELLEFHYYLSSKLNK
LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
NKVKSLRILNTRRKL
MMGMFNMLSTVLGVS
AKSSPAYPSVLGQTI
RHLIFCHSKKKCDELAAK
```

**2. create N random groups**

**g₁**

```
----IDQVTIAGAKLRSLN-
WIQKETLVTFKNPHAKKQDV
--ELLEFHYYLSSKLNK---
---AKSSPAYPSVLGQTI--
```

**g₂**

```
-SLFIGLKGDIRESTV--
---SFSCIAIGIITLYLG
KMLLDNINTPEGIIP---
-LNKYVHGTWRSILP---
--NKVKSLRILNTRRKL-
```

**gₙ**

```
---ESLHNPYPDYHWLRT-
RHLIFCHSKKKCDELAAK-
----VFRLKGGAPIKGVTF
--LNKFISPKSVAGRFA--
DGEEEVQLIAAVPGK----
-MMGMFNMLSTVLGVS---
```

**3** **Simple**          **Remove peptide**

MMGMFNMLSTVLGVS

**5b. Score new core to all log-odds matrices**

**4b. Random shift of the core**

MMGMFNMLSTVLGVS

E₁    E₂    E₃

**6b. Accept or reject move**

$$dE = E_{before} - \max(E_1, E_2, E_3) \qquad P = \min\left[1, \exp\left(\frac{dE}{T}\right)\right]$$

# The scoring function

```
-SLFIGLKGDIRESTV--
---SFSCIAIGIITLYLG
KMLLDNINTPEGIIP---
-LNKYVHGTWRSILP---
--NKVKSLRILNTRRKL-
```

```
-SLFIGLKGDIRESTV--
-SFSCIAIGIITLYLG--
KMLLDNINTPEGIIP---
-LNKYVHGTWRSILP---
--NKVKSLRILNTRRKL-
```

$$LO_{A,j} = \frac{n}{n+\sigma} \log \frac{p_{A,j}'}{q_A}$$

Avoid small specialized clusters ($\sigma$ = 10)

$$E = \sum_j LO_{A,j} \qquad E_i^* = E_i - \lambda \max_{\substack{1 \le n \le g \\ n \ne i}}(E_n, 0)$$

Maximize intra cluster similarity whilst minimize inter cluster similarity

$$dE = E_{before}^* - E_{after}^*$$

$$P = \min\left[1, \exp\left(\frac{dE}{T}\right)\right]$$

P = Probability of accepting the move

**DTU**

Department of Systems Biology
Technical University of Denmark

33

CENTERFO
RBIOLOGI
CALSEQU
ENCEANA
LYSIS CBS

CENTER FOR BIOLOGICAL SEQUENCE ANALYSIS
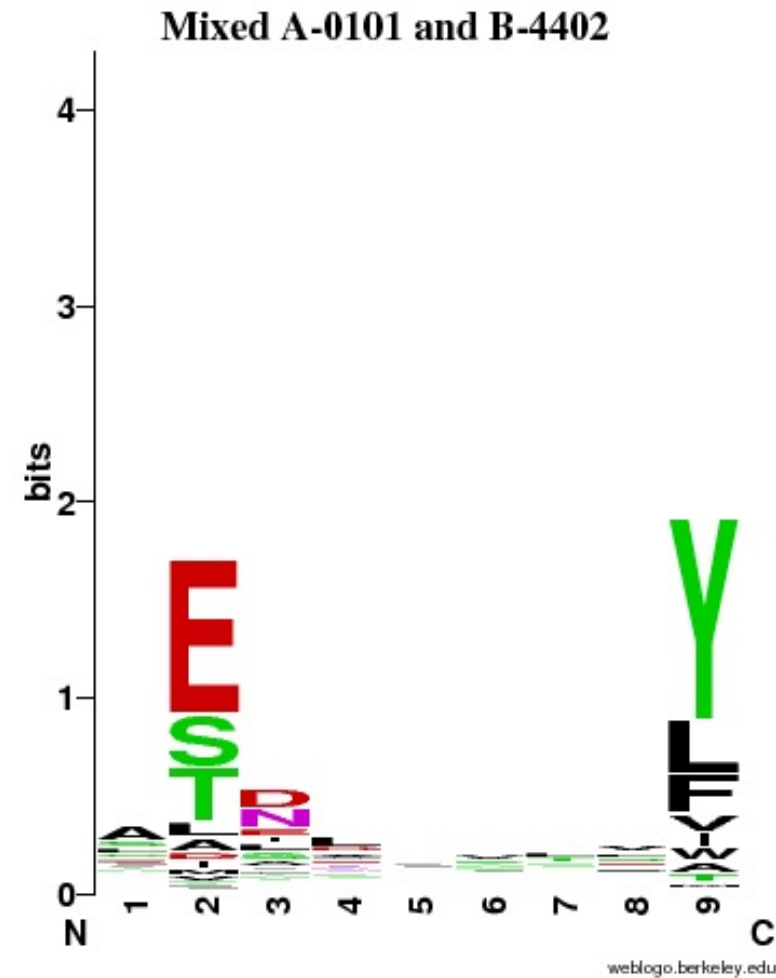
# A mixture of 500 9mer peptides. How many motifs?

# A mixture of 500 9mer peptides. How many motifs?
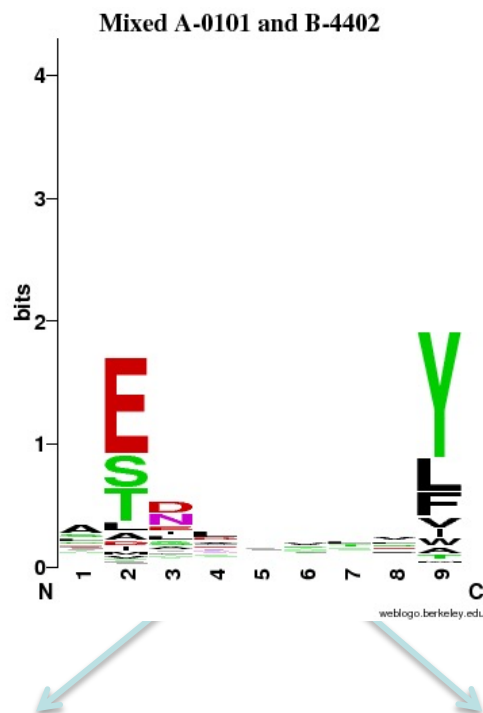
# Two MHC class I alleles: HLA-A*0101 and HLA-B*4402

Mixture of 100 binders
for the two alleles

| Sequence | Allele |
|----------|--------|
| ATDKAAAAY | A*0101 |
| EVDQTKIQY | A*0101 |
| AETGSQGVY | B*4402 |
| ITDITKYLY | A*0101 |
| AEMKTDAAT | B*4402 |
| FEIKSAKKF | B*4402 |
| LSEMLNKEY | A*0101 |
| GELDRWEKI | B*4402 |
| LTDSSTLLV | A*0101 |
| FTIDFKLKY | A*0101 |
| TTTIKPVSY | A*0101 |
| EEKAFSPEV | B*4402 |
| AENLWVPVY | B*4402 |



Mixed A-0101 and B-4402

weblogo.berkeley.edu

DTU

Department of Systems Biology
Technical University of Denmark

Mixed

|     | A0101 | B4402 |
|-----|-------|-------|
| G 1 |       |       |
| G 2 |       |       |

Mixed

Resolved

|      | A0101 | B4402 |
|------|-------|-------|
| G 1  | 97    | 3     |
| G 2  | 3     | 97    |

# Five MHC class I alleles

Mixed 5 MHC I alleles

| | A010 | A020 | A030 | B0702 | B4402 |
|------|------|------|------|-------|-------|
| G 0 | | 1 | 1 | | |
| G 1 | | | | | |
| G 2 | | | | | |
| G 3 | | | | | |
| G 4 | | | | | |

# Five MHC class I alleles

**Mixed 5 MHC I alleles**

|       | A010 | A020 | A030 | B0702 | B4402 |
|-------|------|------|------|-------|-------|
| G 0   | 10   | 11   | 176  | 1     | 0     |
| G 1   | 2    | 4    | 0    | 0     | 95    |
| G 2   | 5    | 87   | 5    | 1     | 0     |
| G 3   | 93   | 2    | 19   | 0     | 2     |
| G 4   | 0    | 6    | 0    | 98    | 3     |



| Group 0 | Group 1 | Group 2 | Group 3 | Group 4 |
|---------|---------|---------|---------|---------|
| HLA-A0301 97% | HLA-A0101 80% | HLA-A0201 89% | HLA-B4402 94% | HLA-B0702 92% |

# Adding in alignment (MHC class II)



HLA-DRB1*03:01    HLA-DRB1*04:01

# Dealing with noisy data

- Experimental data often contain <span style="color:red">false positives</span>

- Outliers do not match any recurrent motif

- Introduce a <span style="color:red">garbage bin</span> to collect outliers

# Dealing with noisy data
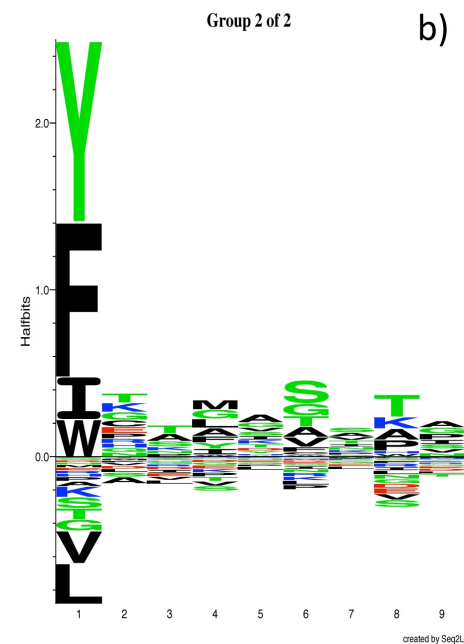
- Introduce a <span style="color:red">garbage bin</span> to collect outliers
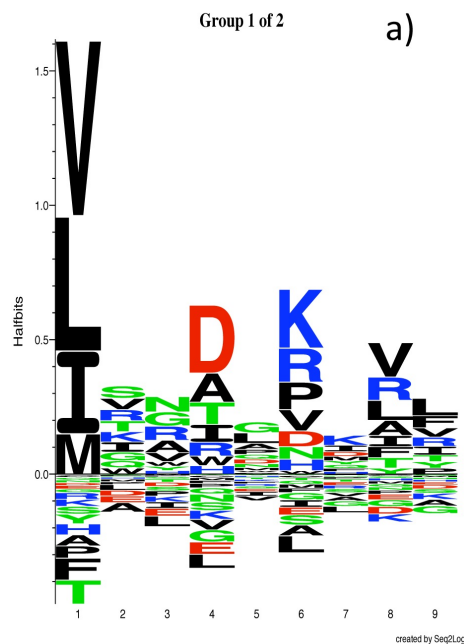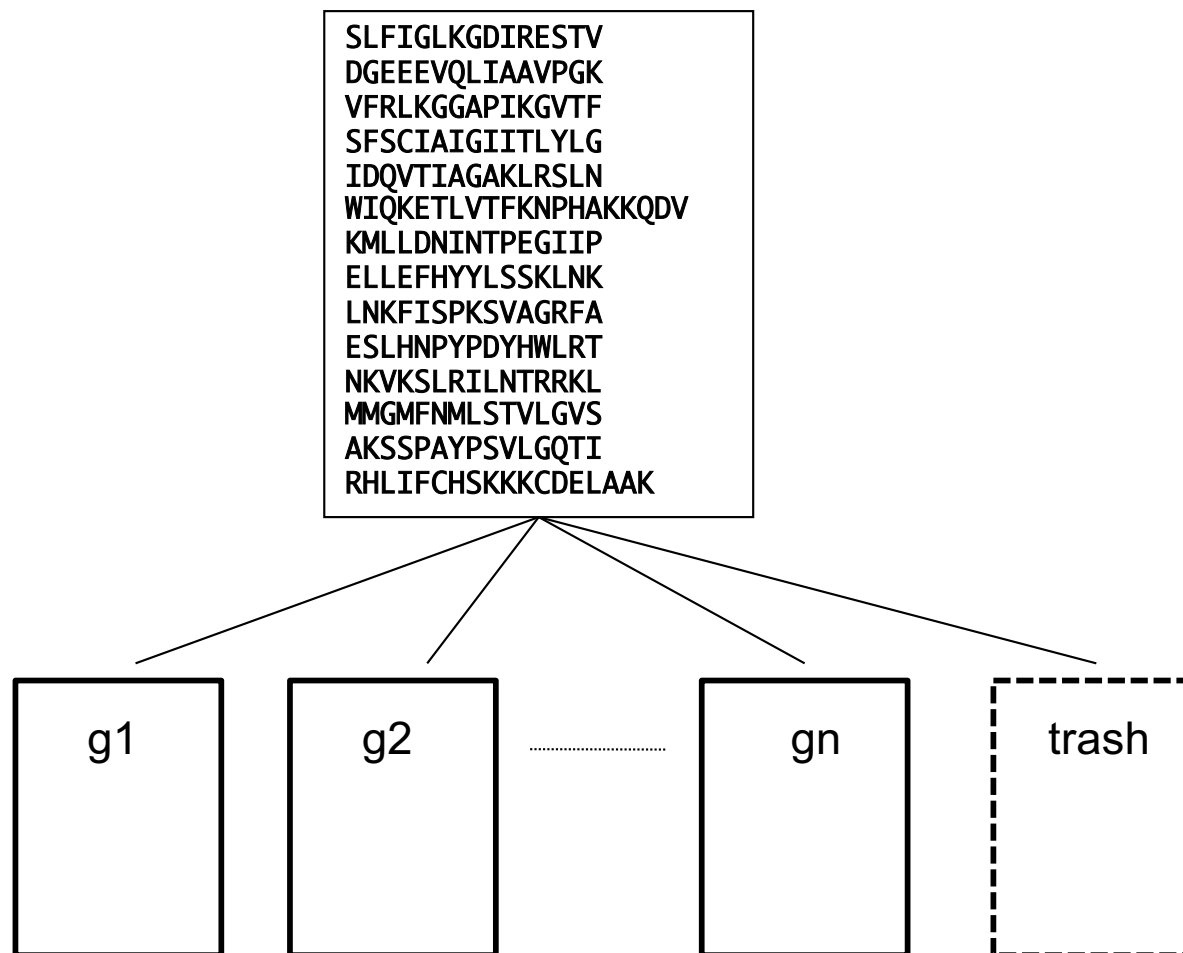
```
SLFIGLKGDIRESTV
DGEEEVQLIAAVPGK
VFRLKGGAPIKGVTF
SFSCIAIGIITLYLG
IDQVTIAGAKLRSLN
WIQKETLVTFKNPHAKKQDV
KMLLDNINTPEGIIP
ELLEFHYYLSSKLNK
LNKFISPKSVAGRFA
ESLHNPYPDYHWLRT
NKVKSLRILNTRRKL
MMGMFNMLSTVLGVS
AKSSPAYPSVLGQTI
RHLIFCHSKKKCDELAAK
```

| g1 | g2 | ..................... | gn | trash |

Move to the
trash cluster
peptides that
do not match
any motif

**DTU**

Department of Systems Biology
Technical University of Denmark

# Dealing with noisy data

## 200 binders to 3 MHC class I alleles

50 random sequences are added to the data set

| 3 alleles | HLA-A0101 | HLA-B0702 | HLA-B4001 | Random | |
|---|---|---|---|---|---|
| g0 | 0 | 197 | 2 | 6 | 205 |
| g1 | 199 | 1 | 0 | 2 | 202 |
| g2 | 0 | 0 | 196 | 5 | 201 |
| trash | 1 | 2 | 2 | 37 | 42 |
| | 200 | 200 | 200 | 50 | |

# Dealing with noisy data

## 200 binders to 3 MHC class I alleles

50 random sequences are added to the data set

| 3 alleles | HLA-A0101 | HLA-B0702 | HLA-B4001 | Random | |
|-----------|-----------|-----------|-----------|--------|-----|
| g0 | 0 | 197 | 2 | 6 | 205 |
| g1 | 199 | 1 | 0 | 2 | 202 |
| g2 | 0 | 0 | 196 | 5 | 201 |
| trash | 1 | 2 | 2 | 37 | 42 |
| | 200 | 200 | 200 | 50 | |

DHHFTPQII

NAFGWENAY
SQTSYQYLI

ELPIVTPAL
ADKNLIKCS

# Dealing with noisy data

**Table 1:** Measured, predicted and re-tested binding affinities (in nM) for peptides assigned to the trash cluster.

| Peptide | HLA | IEDB [a] | Predicted [b] | Validated [c] |
|---------|-----|----------|---------------|---------------|
| DHHFTPQII | A*01:01 | 62 | 28485 | 24822 |
| SQTSYQYLI | B*07:02 | 248 | 24349 | 49928 |
| NAFGWENAY | B*07:02 | 350 | 24481 | - |
| TVFKGFVNK | B*27:05 | 235 | 13723 | - |
| ELPIVTPAL | B*40:01 | 314 | 15208 | - |
| ADKNLIKCS | B*40:01 | 316 | 33324 | 76190 |

[a] Binding affinity deposited in the Immune Epitope Database.
[b] Predicted binding affinities using NetMHCcons.
[c] Re-tested binding affinities after detection as outliers.
As a rule of thumb, generally affinity<50nM identifies a strong binder,
50nM<affinity<500nM a weak binder, affinity>500nM non-binders.

# http://www.cbs.dtu.dk/services/GibbsCluster

CENTERFOR
RBIOLOGI
CALSEQU
ENCEANA
LYSIS CBS

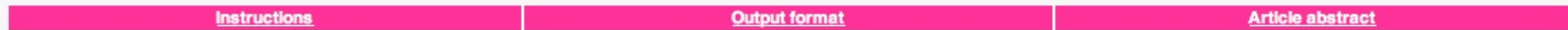| | | | | | | |
|---|---|---|---|---|---|---|
| EVENTS | NEWS | RESEARCH GROUPS | CBS PREDICTION SERVERS | CBS DATA SETS | PUBLICATIONS | EDUCATION |
| STAFF | CONTACT | ABOUT CBS | INTERNAL | CBS BIOINFORMATICS TOOLS | CBS COURSES | OTHER BIOINFORMATICS LINKS |

CBS >> CBS Prediction Servers >> GibbsCluster-1.0

## GibbsCluster-1.0 Server

### Simultaneous alignment and clustering of peptide data

View the version history of this server. All the previous versions are available online, for comparison and reference.

**Instructions**   **Output format**   **Article abstract**

## DATA SUBMISSION

**Paste peptides in the box:**

**or submit a file directly from your local disk:**

Choose File   no file selected

Sample data: Sample 1 - Sample 2

## SUBMIT job

Submit query   Clear fields

DTU

Department of Systems Biology
Technical University of Denmark