

# ALGORITHMS OF BIOINFORMATICS

# 1

## Puzzle from the Lab

*16 October 2025*

Prof. Dr. Sebastian Wild

# Outline

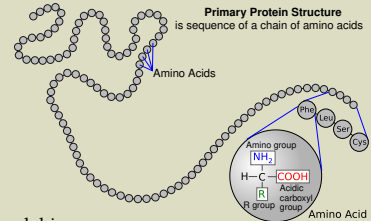
## 1 Puzzle from the Lab

- 1.1 Protein Sequencing
- 1.2 The Turnpike Problem
- 1.3 Backtracking Algorithm
- 1.4 A Pseudopolynomial Algorithm
- 1.5 Back to the Lab

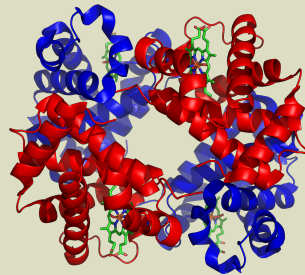
## 1.1 Protein Sequencing

# Proteins: The Workhorses of the Cell

- ▶ **What are they?** Chains of amino acids, folded into specific 3D shapes. The shape determines the function.
- ▶ **What do they do?** Almost everything!
  - ▶ They act as *enzymes* (catalyzing chemical reactions)
  - ▶ provide structural support (cell walls, muscles!),
  - ▶ transport molecules (e. g., *hemoglobin*),
  - ▶ send signals (some *hormones*, e. g., *insulin*)
  - ▶ and more



3D Structure of hemoglobin



[https://commons.wikimedia.org/wiki/File:1GZX\\_Haemoglobin.png](https://commons.wikimedia.org/wiki/File:1GZX_Haemoglobin.png)

⇨ Target of many activities  
across bioinformatics

- ▶ analyzing amino acid sequence
- ▶ predicting structure (AlphaFold)
- ▶ study interaction networks
- ▶ design new proteins  
as potential drugs
- ▶ ...

# Amino Acids

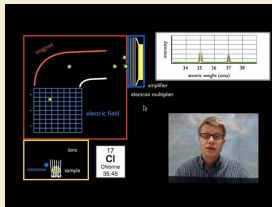
Amino acid	3-letter code	Molecular formula	Mass (Da)
Alanine	Ala	C <sub>3</sub> H <sub>5</sub> NO	71.03711
Cysteine	Cys	C <sub>3</sub> H <sub>5</sub> NOS	103.00919
Aspartic acid	Asp	C <sub>4</sub> H <sub>5</sub> NO <sub>3</sub>	115.02694
Glutamic acid	Glu	C <sub>5</sub> H <sub>7</sub> NO <sub>3</sub>	129.04259
Phenylalanine	Phe	C <sub>9</sub> H <sub>9</sub> NO	147.06841
Glycine	Gly	C <sub>2</sub> H <sub>3</sub> NO	57.02146
Histidine	His	C <sub>6</sub> H <sub>7</sub> N <sub>3</sub> O	137.05891
Isoleucine	Ile	C <sub>6</sub> H <sub>11</sub> NO	113.08406
Lysine	Lys	C <sub>6</sub> H <sub>12</sub> N <sub>2</sub> O	128.09496
Leucine	Leu	C <sub>6</sub> H <sub>11</sub> NO	113.08406
Methionine	Met	C <sub>5</sub> H <sub>9</sub> NOS	131.04049
Asparagine	Asn	C <sub>4</sub> H <sub>6</sub> N <sub>2</sub> O <sub>2</sub>	114.04293
Proline	Pro	C <sub>5</sub> H <sub>7</sub> NO	97.05276
Glutamine	Gln	C <sub>5</sub> H <sub>8</sub> N <sub>2</sub> O	128.05858
Arginine	Arg	C <sub>6</sub> H <sub>12</sub> N <sub>4</sub> O	156.10111
Serine	Ser	C <sub>3</sub> H <sub>5</sub> NO <sub>2</sub>	87.03203
Threonine	Thr	C <sub>4</sub> H <sub>7</sub> NO <sub>2</sub>	101.04768
Valine	Val	C <sub>5</sub> H <sub>9</sub> NO	99.06841
Tryptophan	Trp	C <sub>11</sub> H <sub>10</sub> N <sub>2</sub> O	186.07931
Tyrosine	Tyr	C <sub>9</sub> H <sub>9</sub> NO <sub>2</sub>	163.06333

- ▶ **Dalton (Da):** unit of molecular mass.
- ▶ **1 Da** =  $\frac{1}{12}$  of a carbon-12 atom  
 $\approx 1.66 \times 10^{-27}$  kg.
  - ▶ We will use rounded integer weights
- ▶ **Monoisotopic mass:** sum of atomic masses of most abundant isotopes.
- ▶ Only shows 20 *proteinogenic* amino acids (those encoded in DNA)

# Protein Sequencing

How to determine the sequence of amino acids in a protein?

- ▶ indirect option: via *genes*
  - ▶ ... we will come back to that
  - ▶ not always possible (e. g., for *non-ribosomal peptides*)
- ▶ (more) direct option: *mass spectrometry*
  1. Shatter (many copies) molecule into pieces
  2. Measure *spectrum* of particle masses\* (which masses occur how often)



▶ Mass Spectrometry  
<https://youtu.be/mBT73Pesiog>

⇒ from this, reconstruct what the molecule was!?

## **1.2 The Turnpike Problem**

# Turnpike Problems



▶ The Sopranos Opening

<https://youtu.be/mJpNmYeooQE>

↪ Turnpike = toll road

- ▶ typically, price for road  $\propto$  length of segment on road
- ▶ Can enter and leave at any pair of exits



# Ideal Spectra

Back to mass spectrometry ...

## Simplifying assumptions

- ▶ perfect integer molecular weights, no isotopes
- ▶ all breakpoints realized
- ▶ multiplicities of weights correctly observed
- ▶ no contamination

### Definition 1.1 (Difference multiset)

Given  $P = P[0..n] \in \mathbb{N}_{\geq 1}^n$  a sequence of numbers, define the *prefix sums*  $S[0..n] = \text{prefSum}(P[0..n])$  via  $S[i] = P[0] + \dots + P[i-1]$ .

The *difference multiset*  $\Delta S$  is the multiset

$$\Delta S = \left\{ \left\{ S[j] - S[i] : 0 \leq i < j \leq n \right\} \right\}.$$

Important: Keep duplicates / multiplicities of distances!  $\rightsquigarrow |\Delta S[0..n]| = \binom{n+1}{2}$

# The Turnpike Problem

## Definition 1.2 (Turnpike Problem)

**Given:** a multiset  $D$  with  $|D| = \binom{n}{2}$

**Goal:** Find sequence  $P$  with  $\Delta(\text{prefSum}(P)) = D$  (or state that no such  $P$  exists). ◀

### Examples:

1.  $P_1 = [3, 5, 1, 2]$

$\rightsquigarrow S_1 = [0, 3, 8, 9, 11]$

$\rightsquigarrow D_1 = \Delta S_1 = \{\{1, 2, 3, 3, 5, 6, 8, 8, 9, 11\}\}$

2.  $P_2 = [1, 1, 1, 1, 1]$

$\rightsquigarrow S_2 = [0, 1, 2, 3, 4, 5]$

$\rightsquigarrow D_2 = \Delta S_2 = \{\{1, 1, 1, 1, 1, 2, 2, 2, 2, 3, 3, 3, 4, 4, 5\}\}$

3. For  $D = \{\{1, 1, 1\}\}$  no set  $S$  exists such that  $D = \Delta S$

Any two points  $a < b$  will give  $\Delta(0, a, b) = \{\{a, b, b - a\}\}$  ⚡  $a \neq b$

## 1.3 Backtracking Algorithm

## Systematic Solution

Consider  $\Delta S = \{1, 2, 3, 4, 5, 6, 7, 8, 10, 11, 13, 14, 15, 17, 18\}$ .

# Backtracking Turnpike

---

```
1 procedure turnpikeBacktracking( $D$ )
2    $d := \max D$ 
3    $S := \{0, d\}$  // sorted set of prefSums
4   return turnpikeRec( $S, D$ )
5
6 procedure turnpikeRec( $S, D$ )
7   // Invariant:  $\Delta S \subseteq D$ 
8   if  $\Delta S == D$ 
9     return  $S$ 
10   $d := \max(D \setminus \Delta S)$ 
11  // Option 1: Distance  $d$  from left end
12   $S' := S \cup \{d\}$ 
13  if  $\Delta S' \subseteq D$ 
14     $R := \text{turnpikeRec}(S', D)$ 
15    if  $R \neq \text{NO\_DIFFERENCE\_MULTISET}$ 
16      return  $R$ 
17  // else try Option 2: Distance  $d$  from right
18   $S' := S \cup \{(\max D) - d\}$ 
19  if  $\Delta S' \subseteq D$ 
20    return turnpikeRec( $S', D$ )
21  else // no option worked!
22    return NO_DIFFERENCE_MULTiset
```

---

## ► Correctness

- After placing a few points in prefix sums  $S$ , largest remaining distance must be measured from one endpoint.
- Otherwise we are immediately missing a larger distance ⚡
- ↪ only two checked options are possible
- invariant explicitly checked for recursive calls
- invariant at return guarantees correct answer

## ► Running time

- worst case: exponential! ↪ see tutorials
- not known whether problem is NP-hard(!)

## **1.4 A Pseudopolynomial Algorithm**

# Algebra to the Rescue

Few other algorithmic approaches known for the Turnpike Problem ...  
but one seemingly magic one does!

▶ Consider again  $S = [0, 3, 8, 9, 11]$   $\rightsquigarrow$   $D = \Delta S = \{\{1, 2, 3, 3, 5, 6, 8, 8, 9, 11\}\}$

▶ We can get all pairwise combinations (distances) via *convolutions*

▶ Write  $S(z) = \sum_{s \in S} z^s = z^{11} + z^9 + z^8 + z^3 + z^0$

▶ Now observe that

$$\begin{aligned} S(z) \cdot S(z^{-1}) &= \left( \frac{1}{z^{11}} + \frac{1}{z^9} + \frac{1}{z^8} + \frac{1}{z^3} + 1 \right) (z^{11} + z^9 + z^8 + z^3 + 1) \\ &= z^{11} + z^9 + 2z^8 + z^6 + z^5 + 2z^3 + z^2 + z^1 \\ &\quad + \frac{1}{z^{11}} + \frac{1}{z^9} + \frac{2}{z^8} + \frac{1}{z^6} + \frac{1}{z^5} + \frac{2}{z^3} + \frac{1}{z^2} + \frac{1}{z} + 5 \\ &= \sum_{s \in S} \sum_{t \in S} z^{s-t} \\ &= \sum_{d \in D} z^d + \sum_{d \in D} z^{-d} + |S| \end{aligned}$$

# Factoring Polynomials

- ▶ The expanded product depends only on  $D$ 
  - ↪ can be constructed from the input
- ▶ Use polynomial factorization to check if it can be written as a product  $S(z)S(z^{-1})$ 
  - ▶ this can be done in *pseudopolynomial time*
    - ▶ a polynomial of degree  $d$  with integer coefficients represented with  $b$  bits can be factored over the integers in time  $O(\text{poly}(d, b))$
    - ▶ *Lenstra-Lenstra-Lovász (LLL) algorithm*
  - ▶ polynomial running time in terms of  $n = |D|$ , but exponential in  $b = \log(\max D)$   
 $b$  is the number of bits in the occurring numbers

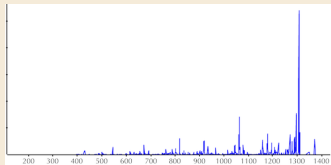


## **1.5 Back to the Lab**

# Ideal vs. Real Spectra

Real protein sequencing tasks unfortunately need additional work . . .

Actual spectrum



Compeau & Pevzner, *Bioinformatics Algorithms*, Fig. 4.13  
<https://cogniterra.org/lesson/29918/step/2?unit=22015>

Values of peaks

372.2	397.2	402.0	<b>406.3</b>	415.1	<b>431.2</b>	<b>448.3</b>	449.3	452.2
471.3	<b>486.3</b>	<b>488.2</b>	500.5	<b>505.3</b>	516.1	536.1	<b>544.2</b>	<b>545.3</b>
562.5	571.3	599.2	614.4	615.4	616.4	618.2	<b>632.0</b>	655.5
656.3	<b>672.5</b>	<b>673.3</b>	677.3	<b>691.4</b>	<b>692.4</b>	712.1	722.3	<b>746.5</b>
760.4	761.6	762.5	<b>771.6</b>	788.4	802.3	803.3	818.5	<b>819.4</b>
831.4	<b>836.3</b>	853.3	875.5	<b>876.5</b>	901.5	915.9	916.5	917.8
<b>918.4</b>	<b>933.4</b>	<b>934.7</b>	<b>935.5</b>	949.4	<b>966.2</b>	995.4	1015.6	1027.5
1029.5	1031.5	1044.5	1046.5	<b>1061.5</b>	<b>1063.4</b>	<b>1079.2</b>	1083.7	
1088.4	1093.5	<b>1096.5</b>	1098.4	1158.5	1159.5	<b>1176.6</b>	1177.7	
1178.6	1192.7	<b>1195.4</b>	1207.5	<b>1210.4</b>	<b>1224.6</b>	1252.5	1270.5	
1271.5	1278.6	1279.6	1295.6	1305.6	1306.5	1307.5	1309.6	

Compeau & Pevzner, *Bioinformatics Algorithms*, Fig 4.14  
<https://cogniterra.org/lesson/29918/step/3?unit=22015>

Ideal Spectrum

0	97	99	113	114	128	128	147	147	163	186	227
241	242	244	260	261	262	283	291	333	340	357	388
389	390	390	405	430	430	447	485	487	503	504	518
543	544	552	575	577	584	631	632	650	651	671	672
690	691	738	745	747	770	778	779	804	818	819	835
837	875	892	892	917	932	932	933	934	965	982	989
1031	1039	1060	1061	1062	1078	1080	1081	1095	1136	1159	1175
1175	1194	1194	1208	1209	1223	1225	1322				

Compeau & Pevzner, *Bioinformatics Algorithms*, Fig 4.7  
<https://cogniterra.org/lesson/29912/step/5?unit=22009>

## Complications:

- ▶ inaccuracy of “weights”
- ▶ weights are actually *mass/charge ratios* (often not so bad)
- ▶ missing/**missed peaks**
- ▶ **false peaks**, e. g., from contamination

# Dealing with Real Spectra

Typical situation in bioinformatics!

- ▶ Inaccuracies in the data
  - ▶ can sometimes be cleaned
  - ▶ or avoided with better lab techniques
  - ▶ or averaged out by producing more repetitions
  - ▶ and/or be worked around by **better algorithms!**
- ▶ For example, we can
  - ▶ Find *best fitting* sequence instead of Yes/No (robust algorithms)
  - ▶ Use further domain knowledge (range of molecular weights of amino acids!)

↪ Must deal with possibilities of incorrect results

- ▶ learn how to judge
- ▶ learn how to communicate shortcomings of methods clearly