# ALGORITHMS
# OF
# BIOINFORMATICS

# *1*

# Puzzle from the Lab

*16 October 2025*

Prof. Dr. Sebastian Wild
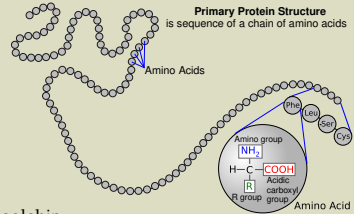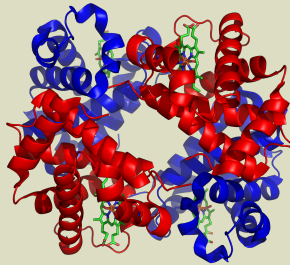
# Outline

## 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



**Primary Protein Structure**
is sequence of a chain of amino acids

Amino Acids

Phe Leu Ser Cys

Amino group
$NH_2$
H — C — COOH
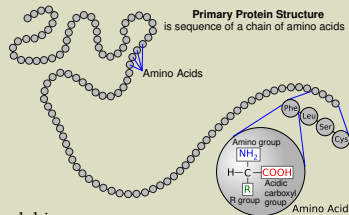Acidic carboxyl group
R group
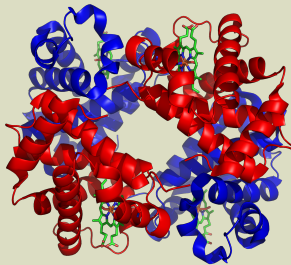Amino Acid

3D Structure of hemoglobin



https://commons.wikimedia.org/wiki/File:1GZX_Haemoglobin.png

# 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

⇝ Target of many activities across bioinformatics
  ► analyzing amino acid sequence
  ► predicting structure (AlphaFold)
  ► study interaction networks
  ► design new proteins as potential drugs
  ► . . .



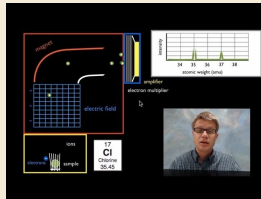**Primary Protein Structure**
is sequence of a chain of amino acids

Amino Acids

Phe
Leu
Ser
Cys

Amino group
$NH_2$
H—C—COOH
Acidic carboxyl group
R group
Amino Acid

3D Structure of hemoglobin



https://commons.wikimedia.org/wiki/File:1GZX_Haemoglobin.png

## Amino Acids

| Amino acid | 3-letter code | Molecular formula | Mass (Da) |
|---|---|---|---|
| Alanine | Ala | $C_3H_5NO$ | 71.03711 |
| Cysteine | Cys | $C_3H_5NOS$ | 103.00919 |
| Aspartic acid | Asp | $C_4H_5NO_3$ | 115.02694 |
| Glutamic acid | Glu | $C_5H_7NO_3$ | 129.04259 |
| Phenylalanine | Phe | $C_9H_9NO$ | 147.06841 |
| Glycine | Gly | $C_2H_3NO$ | 57.02146 |
| Histidine | His | $C_6H_7N_3O$ | 137.05891 |
| Isoleucine | Ile | $C_6H_{11}NO$ | 113.08406 |
| Lysine | Lys | $C_6H_{12}N_2O$ | 128.09496 |
| Leucine | Leu | $C_6H_{11}NO$ | 113.08406 |
| Methionine | Met | $C_5H_9NOS$ | 131.04049 |
| Asparagine | Asn | $C_4H_6N_2O_2$ | 114.04293 |
| Proline | Pro | $C_5H_7NO$ | 97.05276 |
| Glutamine | Gln | $C_5H_8N_2O$ | 128.05858 |
| Arginine | Arg | $C_6H_{12}N_4O$ | 156.10111 |
| Serine | Ser | $C_3H_5NO_2$ | 87.03203 |
| Threonine | Thr | $C_4H_7NO_2$ | 101.04768 |
| Valine | Val | $C_5H_9NO$ | 99.06841 |
| Tryptophan | Trp | $C_{11}H_{10}N_2O$ | 186.07931 |
| Tyrosine | Tyr | $C_9H_9NO_2$ | 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 Problems



▶ The Sopranos Opening
https://youtu.be/mJpNmYeooQE

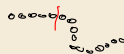$\rightsquigarrow$ Turnpike = toll road

▶ 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

# 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)**   molecular weights / distances / turnpike tolls

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 $S[i] = P[0] + \cdots + P[i-1]$.

## 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)**   <span style="color:brown">molecular weights / distances / turnpike tolls</span>

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 $S[i] = P[0] + \cdots + P[i-1]$.

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! $\rightsquigarrow \left| \Delta S[0..n] \right| = \binom{n+1}{2}$

# 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). ◀

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

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

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

$$③ \; ⑤ \; ① \; ②$$
$$0 \leq i < j \leq n = 4$$
$$\sum_{k=i}^{j-1} P(k)$$

# 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.* $\qquad P_1 = [3, \underline{5}, \underline{1}, 2]$

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

$\quad \leadsto D_1 = \Delta S_1 = \{\!\{1, 2, 3, 3, 5, \underline{6}, 8, 8, 9, 11\}\!\}$

*2.* $\qquad P_2 = [1, 1, 1, 1, 1]$

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

$\quad \leadsto D_2 = \Delta S_2 = \{\!\{1, 1, 1, 1, 1, 2, 2, 2, 2, 3, 3, 3, 4, 4, 5\}\!\}$

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

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

$3\ 5\ 1\ 2$

| | | | |
|---|---|---|---|
| $3$ | $3+5$ | $3+5+1$ | |
| $5$ | $5+1$ | $5+1+2$ | $3+5+1+2$ |
| $1$ | $1+2$ | | |
| $2$ | | | |

$\bigcirc\!\bigcirc$

# Clicker Question

Suppose $\Delta S = \{\{1, 1, 2, 2, 3, 4\}\}$. What is $S$?
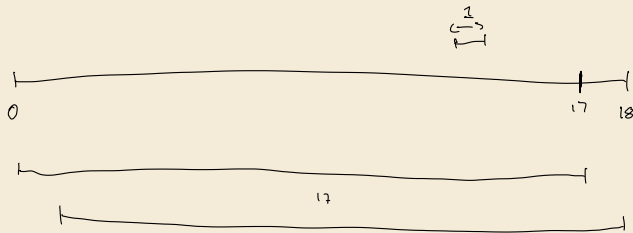
$n = 3$

①  ①  ②

$|$       $|$

$2$  $/$  $|$

→ *sli.do/cs594*

# 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      P := {0, d} // sorted set
4      return turnpikeRec(P, D)
5
6  procedure turnpikeRec(P, D)
7      // Invariant: ΔP ⊆ D
8      if ΔP == D
9          return P
10     d := max(D \ ΔP)
11     // Option 1: Distance d from left end
12     P' := P ∪ {d}
13     if ΔP' ⊆ D
14         R := turnpikeRec(P', D)
15         if R ≠ NO_DIFFERENCE_MULTISET
16             return R
17     // else try Option 2: Distance d from right
18     P' := P ∪ {(max D) − d}
19     if ΔP' ⊆ D
20         return turnpikeRec(P', D)
21     // else: no option worked!
22     return NO_DIFFERENCE_MULTISET
```

# Backtracking Turnpike

```
1  procedure turnpikeBacktracking(D)
2      d := max D
3      P := {0, d} // sorted set
4      return turnpikeRec(P, D)
5
6  procedure turnpikeRec(P, D)
7      // Invariant: ΔP ⊆ D
8      if ΔP == D
9          return P
10     d := max(D \ ΔP)
11     // Option 1: Distance d from left end
12     P' := P ∪ {d}
13     if ΔP' ⊆ D
14         R := turnpikeRec(P', D)
15         if R ≠ NO_DIFFERENCE_MULTISET
16             return R
17     // else try Option 2: Distance d from right
18     P' := P ∪ {(max D) − d}
19     if ΔP' ⊆ D
20         return turnpikeRec(P', D)
21     // else: no option worked!
22     return NO_DIFFERENCE_MULTISET
```

- **Correctness**
    - ▶ After placing a few points in prefix sums $P$, 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

## Backtracking Turnpike

```
1  procedure turnpikeBacktracking(D)
2      d := max D
3      P := {0, d} // sorted set
4      return turnpikeRec(P, D)
5
6  procedure turnpikeRec(P, D)
7      // Invariant: ΔP ⊆ D
8      if ΔP == D
9          return P
10     d := max(D \ ΔP)
11     // Option 1: Distance d from left end
12     P' := P ∪ {d}
13     if ΔP' ⊆ D
14         R := turnpikeRec(P', D)
15         if R ≠ NO_DIFFERENCE_MULTISET
16             return R
17     // else try Option 2: Distance d from right
18     P' := P ∪ {(max D) − d}
19     if ΔP' ⊆ D
20         return turnpikeRec(P', D)
21     // else: no option worked!
22     return NO_DIFFERENCE_MULTISET
```

▶ **Correctness**

   ▶ After placing a few points in prefix sums $P$, 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(!)