# Universitat Politècnica de València

DSIC - Departament de Sistemes Informàtics i Computació

**Academic Project**

MIARFID - Master's in Artificial Intelligence and Pattern
Recognition

# Optimization of Keyboard Layouts for English

**Author:**
Jordi Cantavella Ferrero

Academic Year 2025-2026

# Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

## 1.1  Context

The layout of current keyboards is not optimized for ergonomics and typing efficiency. Current layouts have their roots in traditional typewriting and have not evolved to meet the modern needs of users. Layouts like QWERTY, which emerged in 1873, were designed for mechanical typewriters and are no longer the most efficient for typing on digital devices.

Later studies showed that the QWERTY layout generated inefficient movements, contributing to fatigue and the risk of injury. This led to subsequent layouts such as Dvorak (1936) and Colemak (2006), which aimed to improve efficiency and reduce user fatigue. Seeking to maximize typing speed and minimize physical effort, these layouts are based on ergonomic principles and linguistic analysis.

## 1.2  Objectives

The main objective of this work is the optimization of keyboard layouts, aiming to generate ergonomic and efficient keyboard arrangements. To achieve this, a genetic algorithm and simulated annealing will be implemented to explore the space of possible layouts and find optimal configurations based on ergonomic and linguistic criteria.

- Implement a genetic algorithm for generating and evaluating different

keyboard layouts.

- Implement a simulated annealing algorithm for the same problem.

- Define evaluation metrics that consider distance, hand alternation, and finger effort.

- Analyze the results obtained by the algorithms.

- Compare the generated layouts with existing ones such as *QWERTY*, *Dvorak*, and *Colemak*.

## 1.3   Scope and Limitations

Given the time and resource constraints, this work will focus on optimizing keyboard layouts for the English language. Aspects related to the physical implementation of keyboards or adaptation to different devices will not be addressed.

Also given the different characteristiscs of each distribution, the optimization will be focused on the most common Latin alphabet layouts, including all the letteres, and the most common punctuation marks.

# Chapter 2

# Theoretical Framework

## 2.1 History of Keyboard Layouts

Over time, different types of keyboards have been used, from the first mechanical keyboards to modern digital ones. Among them, the QWERTY, Dvorak, and Colemak keyboards have been particularly prominent.

The **QWERTY** keyboard, designed in 1873 by Christopher Latham Sholes, was created to prevent the keys of early mechanical typewriters from jamming. Its design prioritized separating the most frequently used letters in English to reduce typing speed and prevent the typebars from clashing.

The **Dvorak** keyboard, developed in 1936 by August Dvorak and his brother-in-law William Dealey, was designed to improve efficiency and reduce user fatigue. It demonstrated a reduction of up to 35% in finger movement compared to the QWERTY layout by placing the most used letters on the home row and promoting hand alternation.

The **Colemak** keyboard, introduced in 2006 by Shai Coleman, is an evolution of the QWERTY keyboard that aims to improve efficiency without requiring a radical change in key arrangement. Colemak keeps many keys in their original positions, making the transition easier for users accustomed to QWERTY. It focuses on minimizing finger movement and maximizing typing speed by placing the most common letters in accessible positions.

As shown in Figure 2.1, each of these layouts takes a different approach to

ergonomics and efficiency, reflecting the needs and technologies of its time.
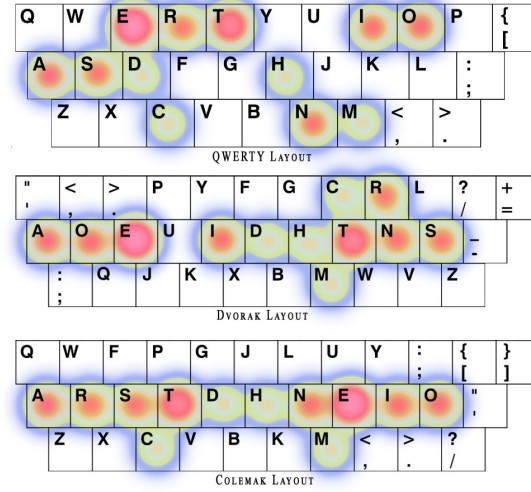


Figure 2.1: Comparison of the QWERTY, Dvorak, and Colemak keyboard layouts.

Representing the different keyboard layouts, Figure 2.1 illustrates how each of these arrangements seeks to optimize the typing experience from different perspectives. Note the concentration of the most used letters on the home row of the Dvorak and Colemak keyboards, in contrast to the more scattered arrangement of the QWERTY keyboard.

## 2.2 Keyboard Ergonomics

Following the main layouts introduced in section 2.1, keyboard ergonomics is a crucial aspect in the design of keyboard layouts. The hand has a series of anatomical characteristics that influence comfort and efficiency when typing. Among them, we highlight:

- **Finger Strength:** The strength of the fingers varies, with the thumb and middle finger being the strongest, while the pinky is the weakest.

- **Finger Mobility:** Fingers have different ranges of motion, with the thumb and index finger being the most mobile.

- **Muscle Fatigue:** Prolonged use of certain fingers can lead to muscle fatigue, especially if constant effort is required.

- **Hand Posture:** A natural and relaxed hand posture is essential to avoid strain and long-term injuries.

## 2.3   Linguistic Analysis

Linguistic analysis is fundamental for designing keyboard layouts that optimize typing efficiency. Paying attention to the frequency of letters and their distribution will allow us to design keyboards that minimize finger movement and maximize typing speed.

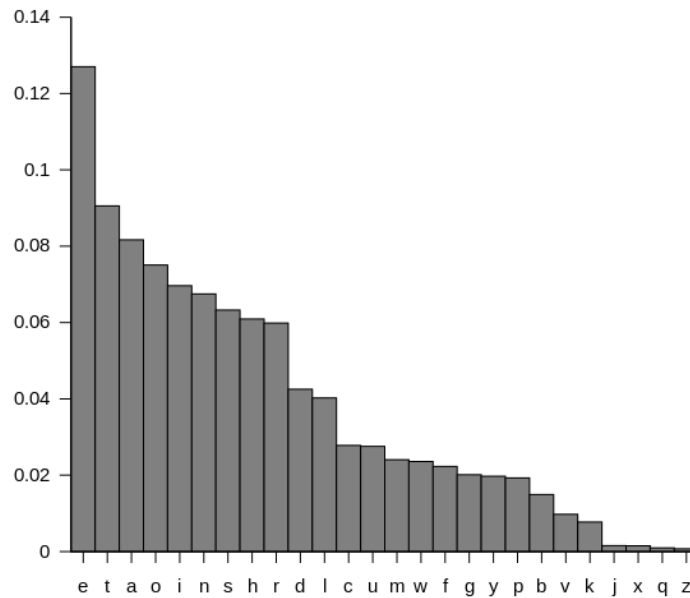Figure 2.2 shows an analysis of letter frequency in English.



Figure 2.2: Frequency of letters in English.

Figure 2.3 shows the frequency distribution of the 26 most common letters used in Latin-based alphabets, covering all languages that use the Latin alphabet, such as Spanish, English, French, German, and Italian. It can be observed that the frequency of letters varies among different languages,

which influences the optimal design of keyboard layouts for each language, even though these languages share a similar alphabet and common origin.
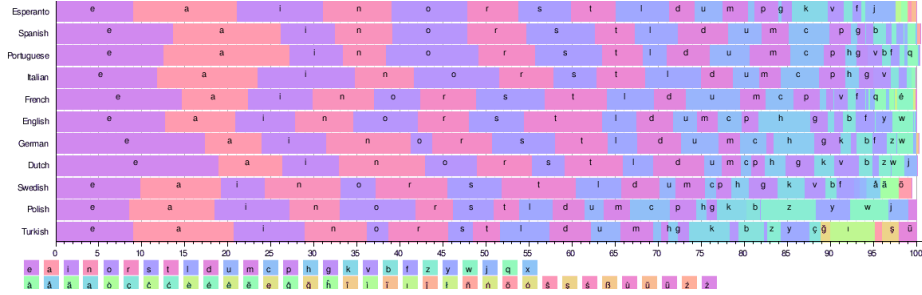


Figure 2.3: Frequency of letters in different languages.

As observed in Figure 2.3, the frequency varies across different languages. It is also important to note that our reference corpus must capture a diverse range of language use, including formal and informal texts, to obtain an accurate representation of the frequency of letters and letter combinations.

## 2.4   Used datasets

In order to analyse the different types of texts and their characteristics, this works uses different types of books or texts as their testing datasets. For a more classical literature approach, we have used *Moby Dick* by Herman Melville, a classic work of English literature. This novel provides a rich and varied vocabulary, making it suitable for analyzing letter and bigram frequencies in traditional literary texts.

Representing a simpler style, we use *The Wonderful Wizard of Oz* by L. Frank Baum. This book features a more constrained and common vocabulary, simpler sentence structures, and repetitive phrasing typical of children's literature. A layout optimized for this corpus should reflect the patterns of more common, everyday English usage.

For each corpus, the raw text will be pre-processed to normalize the data: all letters will be converted to lowercase, and all punctuation except the used in the algorithms, numbers, and whitespace will be removed. From the resulting continuous stream of characters, we will compute the fitness

of each individual in order to determine how good they are as a keyboard layout.

# Chapter 3

# Genetic Algorithm Methodology

This section describes the methodology used to implement a genetic algorithm for optimizing keyboard layouts.

## 3.1 Problem Modeling

The keyboard layout optimization problem can be modeled as a combinatorial optimization problem where the goal is to find the best arrangement of keys on a keyboard. In our particular case, we will focus on the 26 letters of the English alphabet and a few punctuation marks. Trying to minimize the fitness function later introduced in section **??**. Each keyboard layout can be represented as a permutation of the set of characters to be arranged on the keyboard. Our considered set of letters and symbols includes the 26 letters of the English alphabet, along with the space character, period (.), comma (,), semicolon (:) and apostrophe (').

To model our individuals, we use a list of characters where each position corresponds to a specific key on the keyboard. The keyboard is divided into three rows and ten columns, giving us a total of 30 positions that will be filled with the characters from our set. Each key can only be occupied by one character, and each character can only appear once in the layout, ensuring that each individual represents a valid permutation.

The keyboard layout is represented as a standard 3X10 grid structure:

```
 0,  1,  2,  3,  4,  5,  6,  7,  8,  9
10, 11, 12, 13, 14, 15, 16, 17, 18, 19
20, 21, 22, 23, 24, 25, 26, 27, 28, 29
```

Each position is assigned to a specific finger based on standard touch typing conventions, with finger assignments considering hand (left/right), finger type (pinky, ring, middle, index), and relative strength values.

## 3.2    Initial Population

The initial population is created by combining well knowon existing layouts with randomly generated layouts, providing the genetetic a good starting point for exploration. The `init_population` function generates a specified number of individuals, each representing a different keyboard layout, using the following approach: When the `known_distributions` parameter is set to True, the algorithm includes four well-established keyboard layouts as part of the initial population:

- `QWERTY`: The most common layout, designed originally for typewriters

- `Dvorak`: Scientifically designed layout focused on efficiency and reduced finger movement

- `QWERTZ`: German variant of QWERTY with Y and Z positions swapped

- `Colemak`: Modern layout designed to improve upon QWERTY while maintaining some key positions

The remaining individuals in the population are generated randomly to ensure diversity. This individuals are created by shuffling the list of characters to be arranged on the keyboard, ensuring that each layout is a valid permutation of the character set. This hybrid approach of combining known layouts with random enables the genetic algorithm to start with some good solutions while also exploring a wide range of possibilities.

Each individual in the population undergoes validation to ensure it contains exactly 30 characters with no duplicates or missing characters from the defined character set.

## 3.3  Fitness Function

The fitness function evaluates the quality of a keyboard layout by calculating the total "typing cost" when typing a given text. Lower fitness scores indicate better layouts. The function considers multiple factors that affect typing efficiency and ergonomics.

For each consecutive pair of characters in the input text, the algorithm calculates:

1. `Euclidean distance`: the geometric distance between key positions on the 3×10 grid

2. `Finger penalty`: additional costs based on ergonomic factors

The `finger_penalty` function implements a comprehensive penalty system based on:

- `Same Finger Usage:` heavy penalty (1.0) when consecutive characters require the same finger, with additional penalty (2.0) for weak fingers (pinky and ring finger).

- `Same Hand Usage:` moderate penalty (1.0) when both characters are typed with the same hand, promoting hand alternation.

Also additional penalties for the vertical movements:

- `Adjacents rows (distance 1):` penalty of 0.2, with additional 0.15 for weak fingers

- `Two-row jumps (distance 2):` Tpenalty of 0.8, with additional 0.5 for weak fingers

Regarding the usage of the fingers, the weaker fingers are penalised, based on:

- `Pinky finger usage:` +0.15 penalty

- `Ring finger usage:` +0.1 penalty

Finally an additional penalisation for the movements between the same column. Giving an additional penalty of 0.3 for vertical movement within the

same column, with extra penalties for outer columns (0.2 for columns 0 and 9, 0.1 for columns 1 and 8).

So each character pair cost is computed based on:

$$fitness\_function = euclidean\_distance \times \max(1.0 + finger\_penalty, 0.1) \tag{3.1}$$

And the objetive is to minimize said fitness functions.

## 3.4 New Population Generation

### 3.4.1 Selection

The algorithm implements tournament selection as the primary selection mechanism. The `tournament_selection` function works as follows:

1. Groups of k elements are established (where k is the size group). In this particular case having k=3 in order to balance between good and bad fitness individuals.

2. From each group we select the best individual (the one with the lowest fitness function).

3. We perform successive tournaments to select both parments for the crossover operation.

Using the tournament selection over other selection methods we ensure a good balance between selection pressure and diversity maintenance, allowing both high-quality individuals and some lower-quality individuals to participate in reproduction.

### 3.4.2 Crossover

The crossover enables to generate new individuals, given two parents. In this particular case we used a system based on permutations. This operation is designed as follows:

1. **Segment Selection:** randomly select start and end positions to define a crossover segment.

2. `Direct Copy:` copy the selected segment from `parent1` to the child.

3. `Mapping Creation:` create a mapping between characters in the crossover segment of both parents.

4. `Completion:` fill remaining positions using `parent2`, applying the mapping to resolve conflicts and ensure valid permutation, based on a reconstruction function.

The crossover operation is applied with probability `crossover_rate` (typically 0.8). If crossover is not applied, the child is a direct copy of parent1.

A `fix_duplicates` function ensures that any potential duplicates are resolved by replacing them with missing characters, maintaining the permutation property.

### 3.4.3  Mutation

The mutation operator implements a simple swap mutation strategy:

1. `Mutation detection:` each individual has a probability mutation_rate (typically 0.25) of being mutated.

2. `Swap Operation:` if mutation occurs, randomly select two positions and swap their characters.

3. `Permutation Preservation:` the swap operation maintains the valid permutation property

This mutation strategy provides local search capability while ensuring that all individuals remain valid keyboard layouts.

### 3.4.4  Replacement

The algorithm uses an elitist replacement strategy that combines the best individuals from the current generation with newly generated offspring:

1. `Elite preservation:` keep the top 50% of the individuals, particularly elite ones.

2. `Offspring integration:` fill the reamining populaton slots with the offspring from corssover and mutation operations.

This approach ensures that good solutions are not lost while allowing new genetic material to enter the population through reproduction operations. And maintaining a new population the same size as the previous generation.

## 3.5   Results and Analysis

# Chapter 4

# Simulated Annealing Methodology

# Chapter 5

# Conclusions

## 5.1  Summary of the developed work

## 5.2  Achieved objectives

## 5.3  Future work