

FUNCTIONAL PROGRAMMING IN PYTHON FOR LINGUISTS

1. Functional Programming

What is Functional Programming?

- Functional programming (FP) is a **programming paradigm** that focuses on **pure functions, immutability, and avoiding side effects**. Unlike imperative programming, where you define a sequence of commands that modify the program's state, functional programming treats computation as the evaluation of mathematical functions.

Contd.

Key Principles of Functional Programming

1. First-Class Functions

1. Functions can be assigned to variables, passed as arguments, and returned from other functions.
2. This makes code more modular and reusable.

2. Immutability

1. Once a variable is assigned, it **cannot** be changed.
2. Instead of modifying existing data, FP creates new data structures.

Contd.

3. Pure Functions

- A function's output depends only on its input and has **no side effects** (i.e., it does not modify external variables or data).
- This makes code easier to test and debug.

4. Higher-Order Functions

- Functions that take other functions as arguments or return functions.
- Common examples: *map()*, *filter()*, and *reduce()*.

Contd.

5. Recursion Instead of Loops

- Instead of using *for* or *while* loops, FP relies on **recursive function calls**.

6. Lazy Evaluation

- Expressions are not evaluated until their values are needed.
- This improves performance and optimizes memory usage.

Contd.

Why is Functional Programming Useful in Linguistics?

- Linguistics involves processing large amounts of text, such as corpora, dictionaries, and phonological patterns.
- Functional programming is particularly useful because it:

1. Handles Large Datasets Efficiently

- Linguistic data (like corpora or speech datasets) can be processed in parallel using functional constructs like *map()*.

2. Encourages Modular & Reusable Code

- Functional programming allows you to write small, reusable functions.
- Example: A function that converts text to lowercase can be reused in multiple linguistic analyses.

Contd.

3. Avoids Unintended Side Effects

- Linguistic processing often involves **multiple transformations** on text (e.g., tokenization, stemming, lemmatization).
- FP ensures **data remains unchanged**, preventing accidental corruption.

4. Supports Parallel Processing

- Corpus linguistics and sentiment analysis require processing **millions of words**.
- Functional programming supports **parallel execution** (e.g., *multiprocessing* in Python).

Contd.

Example: Imperative vs. Functional Style in Text Processing

Consider a task where we need to **filter out stopwords** from a sentence.

- **Imperative Style (Traditional Approach)**

```
1 stopwords = {"a", "the", "in", "on", "and", "is"}
2 sentence = "This is a simple example"
3
4 words = sentence.split()
5 filtered_words = []
6
7 for word in words:
8     if word not in stopwords:
9         filtered_words.append(word)
10
11 print(filtered_words)
```

```
['This', 'simple', 'example']
```


Contd.

Problems with Imperative Style

- Uses a loop and mutable lists (*filtered_words*).
- Code is longer and harder to debug.

Contd.

- **Functional Style** (Using *filter()* and *lambda*)

```
1 stopwords = {"a", "the", "in", "on", "and", "is"}
2 sentence = "This is a simple example"
3
4 filtered_words = list(filter(lambda w: w not in stopwords, sentence.split()))
5
6 print(filtered_words)
7
```

- Output: `['This', 'simple', 'example']`

➤ Why is this Better?

- **More concise:** Eliminates explicit loops.
- **Immutable:** No unnecessary modifications.
- **More readable & reusable.**

2. Variables, State, and Mutability

Understanding Immutability

- In functional programming, variables don't change once assigned.
- Helps prevent bugs caused by unexpected state changes.

Contd.

- **Example:**

```
1 | x = "hello"  
2 | y = x.upper()  
3 | print(y)    # HELLO  
4 | print(x)    # hello (unchanged)
```

Why Avoid State Changes?

- Reduces side effects.
- Makes debugging easier.

3. Functions as First-Class Objects

- **What Does "First-Class Citizen" Mean?**

- Functions can be stored in variables, passed as arguments, and returned from other functions.

- **Example:**

```
1 def greet(name):  
2     return "Hello, " + name  
3  
4 welcome = greet # Assign function to a variable  
5 print(welcome("Alice")) # Output: Hello, Alice  
~
```

- **Linguistics Application:**

- Defining functions for text transformations.

4. Lambda Functions

What is a Lambda Function?

- A small anonymous function using the *lambda* keyword. Unlike regular functions created with `def`, lambda functions do not have a name and are typically used for short, simple operations.

➤ Syntax: `lambda arguments: expression`

- `lambda` is the keyword.
 - `arguments` are the input parameters (like a normal function).
 - `expression` is the single operation the function performs (must return a value).
-
- **Key Difference from `def` Functions:**
 - Lambda functions are one-liners (cannot contain multiple statements).
 - Implicit return (no need to use `return` keyword).

Contd.

- **Example: Converting Words to Uppercase**
- Using a regular function (def):

```
1 def uppercase(word):  
2     return word.upper()  
3  
4 print(uppercase("phoneme")) # Output: PHONEME
```

- Using a lambda function:

```
1 uppercase = lambda word: word.upper()  
2 print(uppercase("phoneme")) # Output: PHONEME
```

- **Why Use Lambda Functions?**
 - Shortens simple functions.
 - Useful in functions like *map()* and *filter()*.

5. Overt Recursion

What is Recursion?

- A function that **calls itself** until a **base condition** is met.
- A recursive function **breaks down a complex problem into smaller, simpler subproblems**.
- Every recursive function must have a **base case** (a stopping condition) and a **recursive step** (where it calls itself).

Contd.

- **Example: Counting Syllables Recursively**

```
1 def count_vowels(word, count=0):  
2     if word == "":  
3         return count  
4     return count_vowels(word[1:], count + (1 if word[0] in "aeiou" else 0))  
5  
6 print(count_vowels("phonology")) # Output: 3
```

➤ **Why is Recursion Useful in Linguistics?**

- Parsing Nested Structures (e.g., Sentence Trees)
- Morphological Analysis (Breaking Words into Morphemes)
- Phonological Rule Application

6. Comprehensions in Python

- Python comprehensions allow for **concise and efficient** creation of **lists, sets, and dictionaries** in a single line.
- They are faster and more readable than traditional loops.

Contd.

- **Types of Comprehensions**

Comprehension Type	Syntax	Example
List Comprehension	[expression for item in iterable if condition]	[w for w in words if "ing" in w]
Set Comprehension	{expression for item in iterable if condition}	{len(w) for w in words}
Dictionary Comprehension	{key: value for item in iterable if condition}	{w: len(w) for w in words}

Contd.

- **Example: Finding Words with 'ing' Using List Comprehension**

```
1 words = ["running", "jump", "singing", "play"]
2 ing_words = [w for w in words if "ing" in w]
3 print(ing_words) # Output: ['running', 'singing']
```

- **Equivalent Traditional Loop**

```
1 words = ["running", "jump", "singing", "play"]
2 ing_words = []
3 for w in words:
4     if "ing" in w:
5         ing_words.append(w)
6 print(ing_words) # Output: ['running', 'singing']
```

Why use List Comprehension?

- **More concise** than loops.
- **Faster** than appending to lists manually.

Contd.

- Comprehensions use **optimized internal loops**, making them **faster** than traditional loops.

Performance Comparison

```
1 import time
2
3 words = ["morphology", "syntax", "phonetics"] * 1000 # Large dataset
4
5 # Using a Loop
6 start = time.time()
7 word_lengths = {}
8 for w in words:
9     word_lengths[w] = len(w)
10 end = time.time()
11 print(f"Loop Time: {end - start:.6f} sec")
12
13 # Using Dictionary Comprehension
14 start = time.time()
15 word_lengths = {w: len(w) for w in words}
16 end = time.time()
17 print(f"Comprehension Time: {end - start:.6f} sec")
```

Contd.

- **Expected Output** (Varies by System):

```
Loop Time: 0.000420 sec  
Comprehension Time: 0.000308 sec
```

 **Use comprehensions whenever possible** for simple transformations—they make code **faster and cleaner!**

Contd.

- **Nested Comprehensions:**

Comprehensions can also **replace nested loops** for more complex operations.

- **Example: Creating Word Bigrams Using Nested List Comprehension**

```
1 words = ["syntax", "morphology"]  
2 bigrams = [(w[i], w[i+1]) for w in words for i in range(len(w) - 1)]  
3 print(bigrams)  
4 # Output: [('s', 'y'), ('y', 'n'), ('n', 't'), ('t', 'a'), ...]
```

Why use nested comprehensions?

- **Avoids deeply nested loops.**
- **More efficient** in handling multiple iterations.

Contd.

- **Why Use Comprehensions?**

Feature	Traditional Loops	Comprehensions
Code Length	Longer	Shorter
Readability	More verbose	More Concise
Performance	Slower	Faster
Best For	Complex Logic	Simple Transformations



Rule of Thumb:

-  Use **comprehensions** for simple list, set, or dictionary transformations.
-  Use **loops** when code requires multiple statements or complex logic.

7. Vectorized Computation in Python

- Vectorized computation is a technique where **operations are applied to entire arrays (vectors) at once, instead of using loops.**
- **Key Features:**
 - ✓ Eliminates explicit loops, making code more **concise and readable.**
 - ✓ Uses **optimized low-level implementations** for speed (e.g., NumPy in Python).
 - ✓ Ideal for processing **large linguistic datasets** efficiently.

7.1. Using NumPy for Linguistics

Python's built-in lists are **not optimized for numerical operations**. Instead, **NumPy arrays** allow efficient vectorized computations.

- **Example: Token Frequency in a Corpus**

```
1 import numpy as np
2
3 # Word frequencies from a corpus
4 freqs = np.array([10, 50, 30, 60])
5
6 # Normalize frequencies (convert to probabilities)
7 norm_freqs = freqs / np.sum(freqs)
8
9 print(norm_freqs)
```

- **How does it work?**

- `np.array([10, 50, 30, 60])` creates a NumPy array.
- `np.sum(freqs)` calculates the total frequency.
- `freqs / np.sum(freqs)` divides each element by the total → producing normalized frequencies.

No loops required! NumPy performs element-wise division faster than traditional loops.

- **Output:** `[0.06666667 0.33333333 0.2 0.4]`

Contd.

- **Why is Vectorized Computation Useful?**

1. **Speeds Up Large-Scale Linguistic Computations**

- Processing **millions of words** in a **text corpus** is computationally expensive. NumPy **optimizes operations using C and Fortran**, making it significantly faster.

2. **More Efficient Than Traditional Loops**

- **Loop-Based Approach (Slow)**

```
1 freqs = [10, 50, 30, 60]
2 total = sum(freqs)
3 norm_freqs = [f / total for f in freqs] # List comprehension
4 print(norm_freqs)
```

-  **Problem:** Python lists require **explicit looping**, which is **slower** for large datasets.

Contd.

- **NumPy Approach (Fast & Efficient)**

```
1 import numpy as np
2 freqs = np.array([10, 50, 30, 60])
3 norm_freqs = freqs / np.sum(freqs) # Vectorized computation
4 print(norm_freqs)
```

✓ **No loops** → Operations happen at the **hardware level** (C/Fortran), making it **much faster**.

Contd.

Example: Text Processing with NumPy

Let's say we have a list of **word lengths** and we want to find the **mean length**.

- **Loop-Based Approach:**

```
1 word_lengths = [4, 7, 5, 6, 8]
2 mean_length = sum(word_lengths) / len(word_lengths)
3 print(mean_length) # Output: 6.0
```

- **NumPy Approach (Faster & Cleaner):**

```
1 import numpy as np
2 word_lengths = np.array([4, 7, 5, 6, 8])
3 mean_length = np.mean(word_lengths) # Vectorized mean calculation
4 print(mean_length) # Output: 6.0
```



NumPy computes directly without a loop!

Contd.

- **Real-World Linguistic Applications of Vectorized Computation**

Application	How Vectorization Helps
Token Frequency Analysis	Normalizing word counts without loops
Phonetic Feature Processing	Vectorized transformations of speech signals
Corpus Statistics	Fast computation of mean word length, type-token ratio, etc.
TF-IDF Computation	Computing term frequency-inverse document frequency efficiently

Contd.

- **Why Use Vectorized Computation?**

Feature	Loops (Traditional Approach)	Vectorized (NumPy)
Speed	Slower (executes element by element)	Faster (executes at C-level)
Code Simplicity	Requires multiple lines of code	Achieves the same in one line
Memory Efficiency	Stores intermediate results	Optimized for memory usage
Performance	Iterates in Python (slower)	Uses optimized C/Fortran operations
Scalability	Becomes slow with large datasets	Handles large corpora efficiently

8. Iterables, Iterators, and Generators in Python

What is an Iterable?

- An **iterable** is any object in Python that can be **looped over** (e.g., lists, tuples, strings, dictionaries, etc.).
- It contains elements that can be accessed **one by one** but does not keep track of where it is in the iteration.
- **Example: Lists (An Iterable)**

```
1 words = ["syntax", "morphology", "phonetics"]  
2 for w in words:  
3     print(w)
```

- ✓ Works with loops (e.g., for loops)
- ✗ Does not remember position once exhausted

Contd.

What is an Iterator?

- An iterator is an object that **remembers its position** in the sequence while iterating.
 - It must implement `__iter__()` and `__next__()` methods.
 - Once an iterator is exhausted, it **cannot be reset** (unlike iterables).
- **Example: Creating an Iterator from an Iterable**

```
1 words = ["syntax", "morphology", "phonetics"]
2 iterator = iter(words)  # Convert list into an iterator
3
4 print(next(iterator))   # Output: syntax
5 print(next(iterator))   # Output: morphology
6 print(next(iterator))   # Output: phonetics
_
```

✗ Calling `next(iterator)` again will raise an error (`StopIteration`).

Contd.

What is a Generator?

- A **generator** is a special type of iterator that **produces values lazily**, meaning it **yields one value at a time** instead of storing everything in memory.
- ✓ Uses **yield** instead of **return**
- ✓ Automatically remembers its state
- ✓ More memory-efficient than lists

Contd.

- **Example: Generator Function for Streaming Words**

```
1 def word_stream():
2     words = ["syntax", "morphology", "phonetics"]
3     for w in words:
4         yield w # Produces one word at a time
5
6 stream = word_stream() # Create generator object
7
8 print(next(stream)) # Output: syntax
9 print(next(stream)) # Output: morphology
10 print(next(stream)) # Output: phonetics
```

➤ **How Does It Work?**

Calling `word_stream()` **does not execute the function immediately** → It creates a **generator object**.

The first `next(stream)` **executes until it reaches yield**, returning "syntax".

The second `next(stream)` **resumes execution** where it left off, returning "morphology".

Contd.

- **Why Use Generators?**

- ✓ **Memory-Efficient for Large Text Processing**

- Instead of **loading an entire corpus into memory**, a generator **fetches one word at a time**.

- ✓ **Faster Performance**

- No need to **precompute and store** all values → **Improves speed** for big data.

- ✓ **Can be used in Streaming Data Processing**

- Perfect for handling **large files, corpora, and databases** without **high RAM usage**.

Contd.

- **Example: Reading a Large Text File Line by Line (Efficiently)**

```
1 def read_large_file(filename):  
2     with open(filename, "r", encoding="utf-8") as file:  
3         for line in file:  
4             yield line.strip() # Yield one line at a time  
5  
6 lines = read_large_file("corpus.txt")  
7  
8 print(next(lines)) # Prints the first line of the file  
9 print(next(lines)) # Prints the second line  
10
```

 Does not load the entire file into memory!

Contd.

- **Why Use Iterators & Generators?**

Feature	Iterable (List, String)	Iterator (iter())	Generator (yield)
Loop Over Data	✓ Yes	✓ Yes	✓ Yes
Remembers Position	✗ No	✓ Yes	✓ Yes
Exhaustible?	✗ No	✓ Yes	✓ Yes
Memory Efficiency	✗ Loads all data	✗ Stores all data	✓ Generates on demand
Best For	Small data sets	Keeping track of iteration	Large-scale text processing

Contd.

Key Takeaways

- **Use Iterables** (list, tuple, string) when working with **small datasets**.
- **Use Iterators** when you need to **manually control iteration**.
- **Use Generators** when working with **large text files, corpora, or real-time data streaming**.

9. Parallel Programming in Python

- Parallel processing **speeds up tasks** by executing multiple operations **simultaneously** instead of one after another (sequential processing).
- **Key Benefits:**
 - **Faster execution** → Speeds up **corpus processing, NLP tasks, and text analysis**.
 - **Utilizes multiple CPU cores** → More efficient for **large datasets**.
 - **Ideal for repetitive linguistic computations** (e.g., phoneme classification, token frequency calculations).

Contd.

Using multiprocessing to Process Words in Parallel

- Python's multiprocessing module **allows multiple CPU cores to process data concurrently**. Below, we convert words **to uppercase in parallel**, using two CPU cores.

```
1 from multiprocessing import Pool
2
3 def process_word(word):
4     return word.upper() # Convert to uppercase
5
6 words = ["phoneme", "syllable", "morpheme"]
7
8 # Create a pool of 2 worker processes
9 with Pool(2) as p:
10     results = p.map(process_word, words) # Process words in parallel
11
12 print(results) # Output: ['PHONEME', 'SYLLABLE', 'MORPHEME']
```

- How Does It Work?**

Pool(2) → Creates a **pool of 2 worker processes**.

map(process_word, words) → **Splits the list** across processes and applies process_word().

Each process runs independently, utilizing multiple CPU cores.

Contd.

- **Why Use multiprocessing Instead of a Loop?**

Feature	Loop (Sequential)	Multiprocessing (Parallel)
Speed	Slower (one task at a time)	Faster (multiple tasks at once)
CPU Usage	Uses one core	Uses multiple cores
Efficiency	Good for small tasks	Best for large datasets
Use Case	Simple word processing	Large-scale corpus analysis

Contd.

Real-World Use Cases in Linguistics

- Fast Corpus Processing (e.g., Tokenizing a Large Text Corpus)

```
1 from multiprocessing import Pool
2 import nltk
3
4 nltk.download("punkt")
5 from nltk.tokenize import word_tokenize
6
7 def tokenize_sentence(sentence):
8     return word_tokenize(sentence)
9
10 sentences = [
11     "Phonetics studies the sounds of speech.",
12     "Morphology examines word structure.",
13     "Syntax analyzes sentence structure."
14 ]
15
16 # Process sentences in parallel
17 with Pool(3) as p:
18     tokenized_sentences = p.map(tokenize_sentence, sentences)
19
20 print(tokenized_sentences)
```

- ✓ Speeds up tokenization for large corpora!

Contd.

Parallel Phoneme Classification

- Instead of looping over **millions of phonemes**, we can **process them in parallel**.

```
1 from multiprocessing import Pool
2
3 phonemes = ["a", "e", "i", "o", "u"]
4
5 def classify_phoneme(phoneme):
6     vowels = {"a", "e", "i", "o", "u"}
7     return (phoneme, "vowel" if phoneme in vowels else "consonant")
8
9 with Pool(2) as p:
10     classified_phonemes = p.map(classify_phoneme, phonemes)
11
12 print(classified_phonemes)
13 # Output: [('a', 'vowel'), ('e', 'vowel'), ('i', 'vowel'), ('o', 'vowel'), ('u', 'vowel')]
```

✓ Handles large phoneme datasets efficiently!

Contd.

When NOT to Use Parallel Processing

✗ If a task **doesn't require high computation**, parallelization may **add overhead** instead of speeding things up.

✓ Use parallel processing for **CPU-intensive tasks**, like **text analysis, speech recognition, and NLP**.

Contd.

- **Final Takeaways:**

- Parallel processing **speeds up** large-scale linguistic tasks.
- Use **multiprocessing.Pool** for CPU-bound tasks.
- Ideal for **tokenization, phoneme classification, and corpus analysis**.

10. Making Nonsense Items Again

Nonsense items are **artificially generated words** that **follow linguistic rules** but do **not exist in the language**.

Why Create Nonsense Words?

- Used in **psycholinguistic experiments** to test **word recognition**.
- Helps in studying **phonotactic constraints** (rules about which sounds can appear together).
- Useful for **modeling language acquisition** (how speakers process new words).

Contd.

Using Functional Programming to Generate Fake Words

- We can use **functional programming** (lambda, random.choice()) to generate nonsense words.

```
1 import random
2
3 # Define possible consonants and vowels
4 consonants = "ptkbgdgm"
5 vowels = "aeiou"
6
7 # Generate a nonsense word (CVC pattern)
8 generate_word = lambda: random.choice(consonants) + random.choice(vowels) + random.choice(consonants)
9
10 # Generate and print a random word
11 print(generate_word()) # Example Output: "tam"
```

How Does It Work?

- `random.choice(consonants)` picks a random consonant.
- `random.choice(vowels)` picks a random vowel.
- `random.choice(consonants)` picks another consonant → forming a **CVC (Consonant-Vowel-Consonant)** structure.

Contd.

Linguistics Application: Phonotactic Constraints

- Phonotactic constraints determine **which sound sequences are allowed** in a language.
- **Example: Generating Nonsense Words for English**
 - Some consonant clusters, like "bn" or "zr", are **not allowed** in English.
We can filter **only valid English-like sequences**:

Contd.

```
1 import random
2
3 # Define consonants and vowels with English phonotactic constraints
4 onset = "ptkbgdgm" # Allowed initial consonants
5 vowel = "aeiou"
6 coda = "tksmn" # Allowed final consonants
7
8 # Generate nonsense words following English-like patterns
9 generate_word = lambda: random.choice(onset) + random.choice(vowel) + random.choice(coda)
10
11 print(generate_word()) # Example Output: "tam"
```



Ensures phonotactically valid words!

Contd.

Expanding to Complex Word Structures

- We can generate longer nonsense words (e.g., CCVCVCC patterns):

```
1 import random
2
3 # Consonant clusters
4 clusters = ["bl", "tr", "st", "gr", "pl"]
5 vowels = "aeiou"
6
7 # Generate a nonsense word (CCVCVCC pattern)
8 generate_complex_word = lambda: random.choice(clusters) + random.choice(vowels) + random.choice(clusters)
9
10 print(generate_complex_word()) # Example Output: "plotr"
```

- Why Is This Useful?
 - ✓ Simulates **natural language patterns**
 - ✓ Used in **linguistic experiments**
 - ✓ Can be **customized for different languages**

Contd.

Why Generate Nonsense Words?

Application	How Does It Help?
Psycholinguistics	Test word recognition & learning
Phonology	Study phonotactic rules
Speech Disorders	Assess pronunciation skills
Computational Linguistics	Generate test data for NLP models

11. Summary

Key Takeaways

- Functional programming avoids side effects and uses immutability.
- Recursion and comprehensions simplify code.
- Parallelism speeds up linguistic data processing.