# Radboud University - Nikhef

High Energy Physics department - ATLAS team

# ChatGPT
# as a Coding Assistant

dr. ir. Uraz Odyurt, dr. Geoffrey Gilles
2023-05-15

Say ChatGPT one more time …!

# Stating the obvious …

- What is ChatGPT?

  **OpenAI**

  | CT | What is ChatGPT? |

  ChatGPT is a large language model developed by OpenAI, based on the GPT-3.5 architecture. It is designed to generate human-like responses to natural language input, allowing users to have conversations with the model as if it were a human. ChatGPT has

- Language model => Any use-case relevant to languages …
  => Writing
  => Summarising
  => Translation
  => Comparison
  => …
  => **Code generation and tasks related to coding**

  | CT | Can you teach me how to cook beef bourguignon? |

  | CT | Can you tell me the arguments in favour of god's existence? |

  | CT | Can you compare the economy of Japan with South Korea? |

# A language model you say …

- We cannot have all possible sentences, but the building blocks …

- "Sentence" as a unit is not granular enough.
  => Let's consider words as **atomic units** …

| I | saw | a | raven | flying | over | the | station |
|---|-----|---|-------|--------|------|-----|---------|

$$P(B \cap A) = P(A) \cdot P(B|A)$$

- The joint probability of B and A occurring means the probability of A occurring, multiplied by the probability of B occurring given that **A has occurred** *(context)*
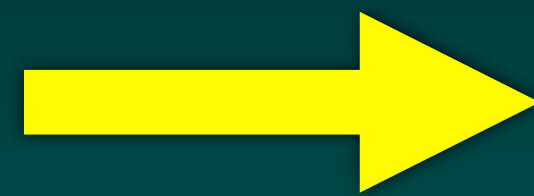
# Extensive knowledge of programming languages

- Ability to understand and interpret natural language queries related to programming
  => Syntax
  => Semantics

- Benefitting from the OpenAI "~~Codex~~" (deprecated -> GPT-3.5)
  => Demo: <u>Creating a Space Game with OpenAI Codex</u>
  => Separate OpenAI **language model**, designed for code generation
  => Interfaced to a large **curated code database**
  => Millions of **code snippets**
  => **Organised** by language, libraries and frameworks
  => Accompanied by **metadata** on function, inputs and outputs

Example => "GitHub Copilot uses the OpenAI Codex to suggest code and entire functions in real-time."

# As a coding assistant

| Error resolution | Code review & optimisation | Algorithm design building blocks | Documentation |
|---|---|---|---|
| • Identifying and resolving coding errors<br>• Guidance - How to fix | • Code analysis and feedback<br>• Optimise code for performance | • Right algorithms<br>• Best libraries and frameworks | • Detailed documentation<br>• Verbose explanations<br>• Simple diagrams |
| • But inability to solve too complex or overarching issues | • Limited by practical experience and context awareness of the user | • Unable to write extensive code, but rather provide relevant code snippets | • Analytical explanations are generic and follow known knowledge/best practices |

# Can we trust the output?



Describe 007?

Johnny English
at your service!

# Can we trust the output?

- Short answer: No!

- **Validation** is the key
  => You need partial/total **expertise** in the relevant field

- Proper answer: It depends …
  => For languages with an **extensive training corpus**: Yes!
  (also programming languages)
  => For **common tasks**: Yes!
  (to save time)

Towards a Benchmark for Scientific Understanding in Humans and Machines

# Let's try it out … !

**CT** Can you explain with simple examples how pytorch view interacts with tensors?

Sure, I'd be happy to explain how PyTorch's `view` function interacts with tensors using some simple examples.

The `view` function in PyTorch is used to reshape a tensor without changing its underlying data. It returns a new tensor with the same data but with a different shape. Here are some examples to demonstrate its usage:

```python
import torch

# create a 1D tensor
x = torch.tensor([1, 2, 3, 4, 5, 6])
print(x)
# Output: tensor([1, 2, 3, 4, 5, 6])

# reshape the tensor into a 2D tensor with 2 rows and 3 columns
y = x.view(2, 3)
print(y)
# Output: tensor([[1, 2, 3],
#                 [4, 5, 6]])

# reshape the tensor into a 3D tensor with 2 channels, 1 row, and 3 columns
z = x.view(2, 1, 3)
print(z)
# Output: tensor([[[1, 2, 3]],
#                 [[4, 5, 6]]])
```

**CT** Can you write the generated data directly to the CSV file, not using a Data Frame?

Yes, you can write the generated data directly to a CSV file without using a DataFrame. Here's an example code that generates a list of tuples containing the same data as in the previous example, and writes it directly to a CSV file:

```python
import csv
import numpy as np

# Generate dummy data
names = ['Alice', 'Bob', 'Charlie', 'David', 'Emily']
ages = np.random.randint(low=20, high=50, size=len(names))
cities = ['New York', 'Los Angeles', 'Chicago', 'Houston', 'Miami']

# Write data to CSV file
with open('dummy_data.csv', mode='w', newline='') as file:
    writer = csv.writer(file)
    writer.writerow(['Name', 'Age', 'City'])
    for i in range(len(names)):
        writer.writerow([names[i], ages[i], cities[i]])
```

**CT** Tell me which sorting algorithm is the fastest and provide the C++ implementation of it?

The fastest sorting algorithm depends on the specific use case and the characteristics of the data being sorted. Generally, quicksort is considered one of the fastest sorting algorithms for most use cases, with an average time complexity of O(n log n). However, it has a worst-case time complexity of O(n^2), which can occur when the input is already sorted or nearly sorted. Other sorting algorithms such as merge sort and heapsort also have an average time complexity of O(n log n), but their worst-case time complexity is better than quicksort.

Here is a C++ implementation of quicksort:

```cpp
#include <iostream>

class Car {
    public:
        // Constructor
        Car(std::string m
        
        // Member functic
        void start() {
            std::cout << "Starting the " << make_ << " "
        }
        void stop() {
            std::cout << "Stopping the " << make_ << " "
        }
    private:
        // Data members
        std::string make_;
        std::string model_;
        int year_;
};

class Lorry : public Car {
```

**CT** Now rewrite the sort in Python.

Here's an implementation of quicksort in Python:

```python
def quicksort(arr):
    if len(arr) <= 1:
        return arr
    else:
        pivot = arr[len(arr)//2]
        left = [x for x in arr if x < pivot]
        middle = [x for x in arr if x == pivot]
        right = [x for x in arr if x > pivot]
        return quicksort(left) + middle + quicksort(right)

arr = [5, 2, 9, 3, 7, 4, 8, 1, 6]
sorted_arr = quicksort(arr)
print(sorted_arr)
```

This implementation is a bit different from the C++ implementation. Instead of using indices to partition the array, it uses list comprehension to create three separate lists for elements less than, equal to, and greater than the pivot element. These lists are then concatenated together in the correct order to produce the sorted array.

The function takes in a list `arr` and recursively calls itself on the left and right partitions until the entire list is sorted. The base case for the recursion is when the length of the list is less than or equal to 1, indicating that it contains only one element.

query and receives back the resulting publications' metadata:

```python
import requests
import json

# Set the API key and search query parameters
api_key = 'YOUR_API_KEY'
query = 'TITLE-ABS-KEY("data mining") AND PUBYEAR AFT 2018'

# Define the API endpoint and request parameters
url = 'https://api.elsevier.com/content/search/scopus'
params = {'apiKey': api_key, 'query': query, 'field': 'dc:identifie

# Send the request to the Scopus Search API and get the response
response = requests.get(url, params=params)

# Parse the response and extract the publications' metadata
if response.status_code == 200:
    data = json.loads(response.text)
    for item in data['search-results']['entry']:
        identifier = item['dc:identifier']
        title = item['dc:title']
        journal = item['prism:publicationName']
        date = item['prism:coverDate']
```

Regenerate response

```python
process = MG.MadGraph5Interface(proc_card='p p > h > gamma gamma',)

# Set the process energy and number of events to generate
process.run_cfg['run_mode'] = 0    # "Run" mode for event generation
process.run_cfg['nb_event'] = 1000  # Number of events to generate
process.run_cfg['energy'] = 13000   # Center-of-mass energy in GeV

# Generate events and save them to a file
process.run_cmd()
```

**CT** Considering the above two snippets, which one has better performance?

Generally, writing data directly to a CSV file using the `csv` module can be more performant than creating a DataFrame and then writing it to a CSV file using pandas. This is because pandas is a powerful and flexible library that provides a lot of functionality, but that comes at a cost of performance overhead.

However, the performance difference between the two methods may not be significant for small datasets. In fact, for small datasets, the overhead of using pandas may be negligible, and the simpler and more concise syntax of using pandas may be more appealing.

In summary, if you're dealing with a large dataset, writing directly to a CSV file using the `csv` module may provide better performance. On the other hand, if you're working with a small dataset, the performance difference may not be significant, and you may prefer to use the simpler and more concise syntax of pandas.