# TYPHON

## Automatic Recommendation of Relevant Code Cells in Jupyter Notebooks

**Chaiyong Ragkhitwetsagul, Veerakit Prasertpol, Natanon Ritta, Paphon Sae-Wong, Thanapon Noraset and Morakot Choetkiertikul**
**Faculty of ICT, Mahidol University**

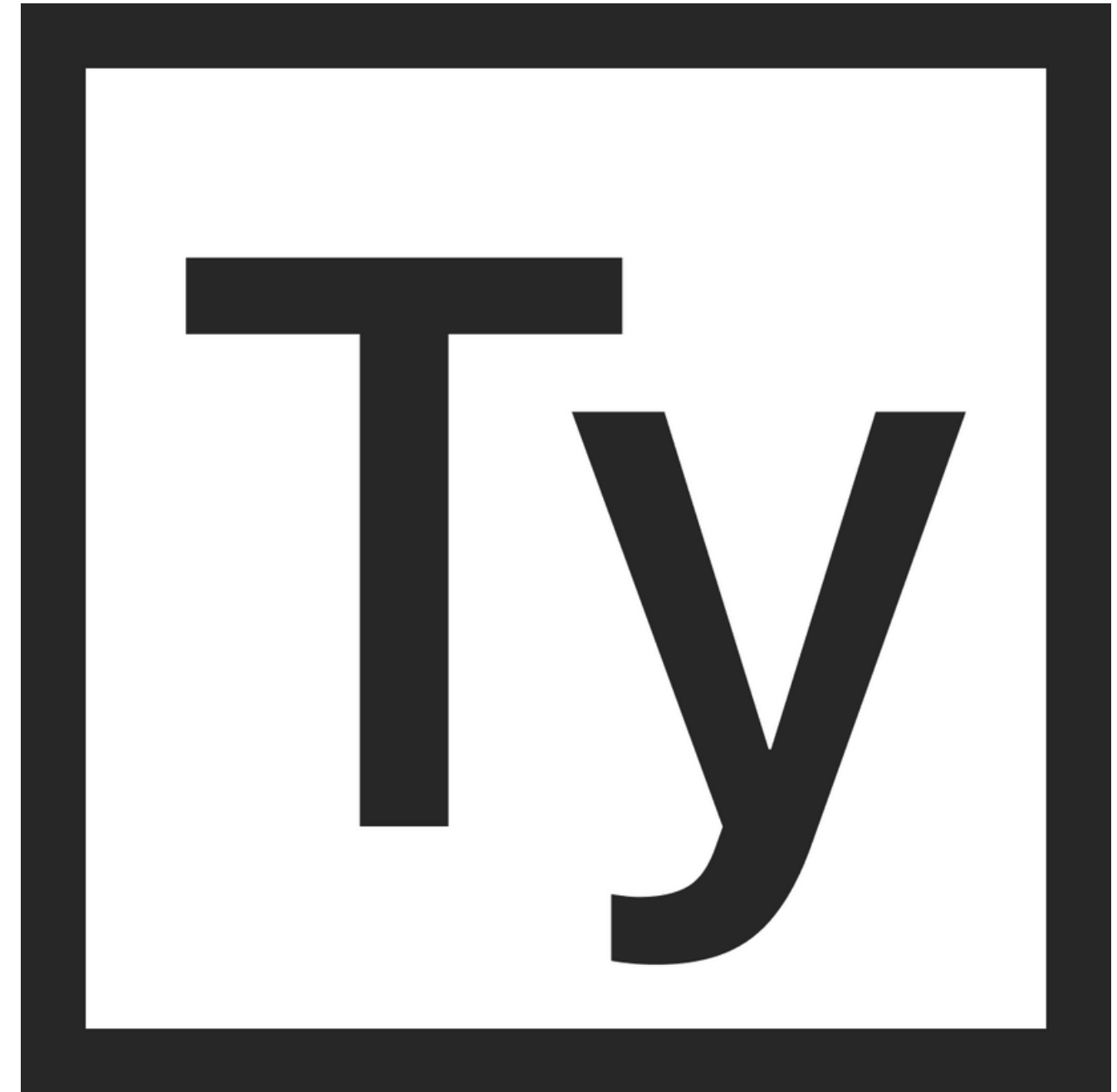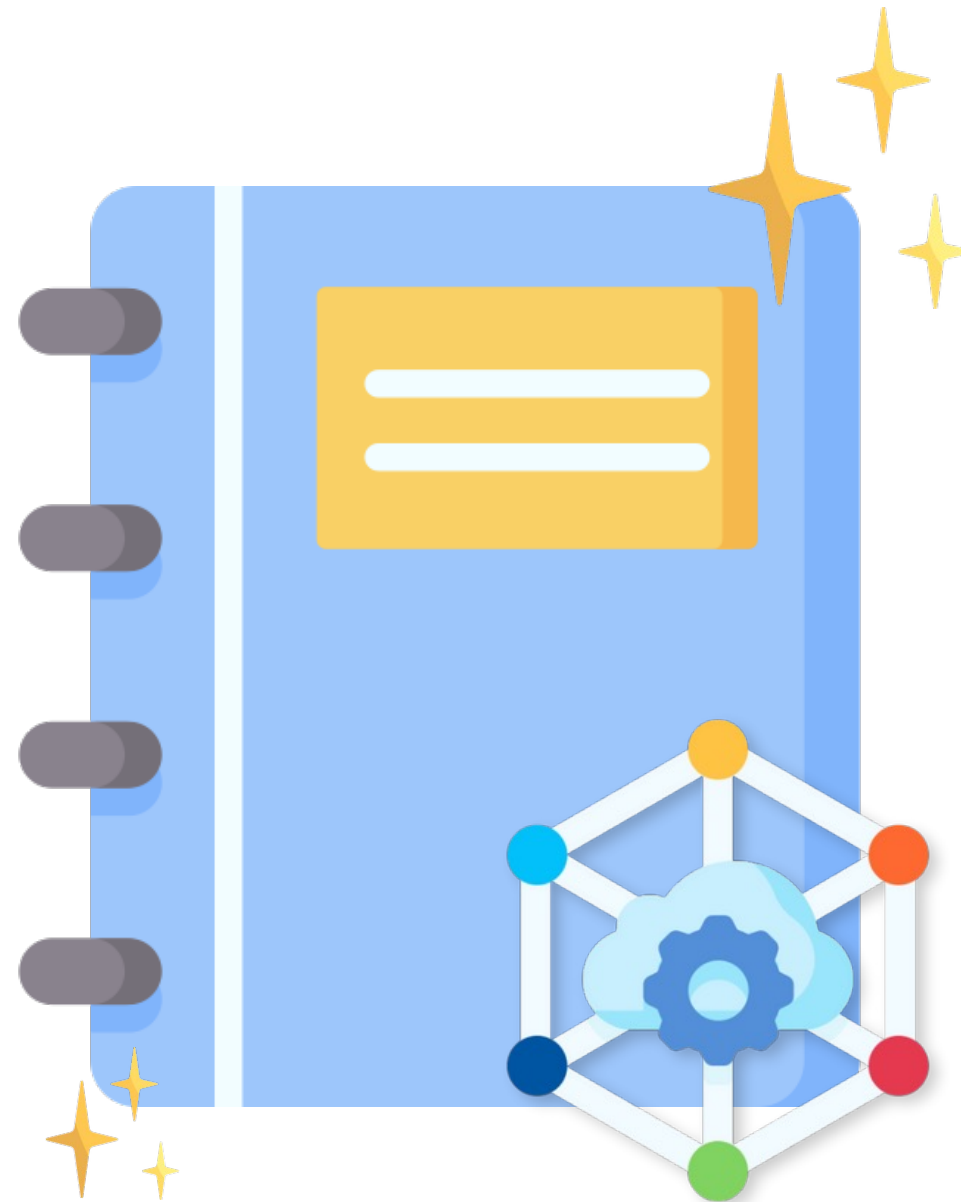Mahidol University
Faculty of Information and Communication Technology

ICT

SERU

# COMPUTATIONAL NOTEBOOK

Computational notebook is a well-known and well-adopted technology in tasks related to data analysis

A Jupyter notebook can be a central place for collaborative data analysis.

File   Edit   View   Insert   Runtime   Tools   Help

+ Code   + Text        Copy to Drive                                                          Connect ▾   ◆ Gemini   ⌃

# Plot styles

Colaboratory charts use Seaborn's custom styling by default. To customize styling further please see the matplotlib docs.

∨ **3D Graphs**

∨ **3D Scatter Plots**

```python
import matplotlib.pyplot as plt
import numpy as np
from mpl_toolkits.mplot3d import axes3d

fig = plt.figure()
ax = fig.add_subplot(111, projection = '3d')

x1 = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
y1 = np.random.randint(10, size=10)
z1 = np.random.randint(10, size=10)

x2 = [-1, -2, -3, -4, -5, -6, -7, -8, -9, -10]
y2 = np.random.randint(-10, 0, size=10)
z2 = np.random.randint(10, size=10)

ax.scatter(x1, y1, z1, c='b', marker='o', label='blue')
ax.scatter(x2, y2, z2, c='g', marker='D', label='green')

ax.set_xlabel('x axis')
ax.set_ylabel('y axis')
ax.set_zlabel('z axis')
plt.title("3D Scatter Plot Example")
plt.legend()
plt.tight_layout()
```

# kaggle

**A CLOUD-BASED COLLABORATIVE PLATFORM INVOLVING DATA ANALYTICS TASKS USING A COMPUTATIONAL NOTEBOOK IN PRACTICES**
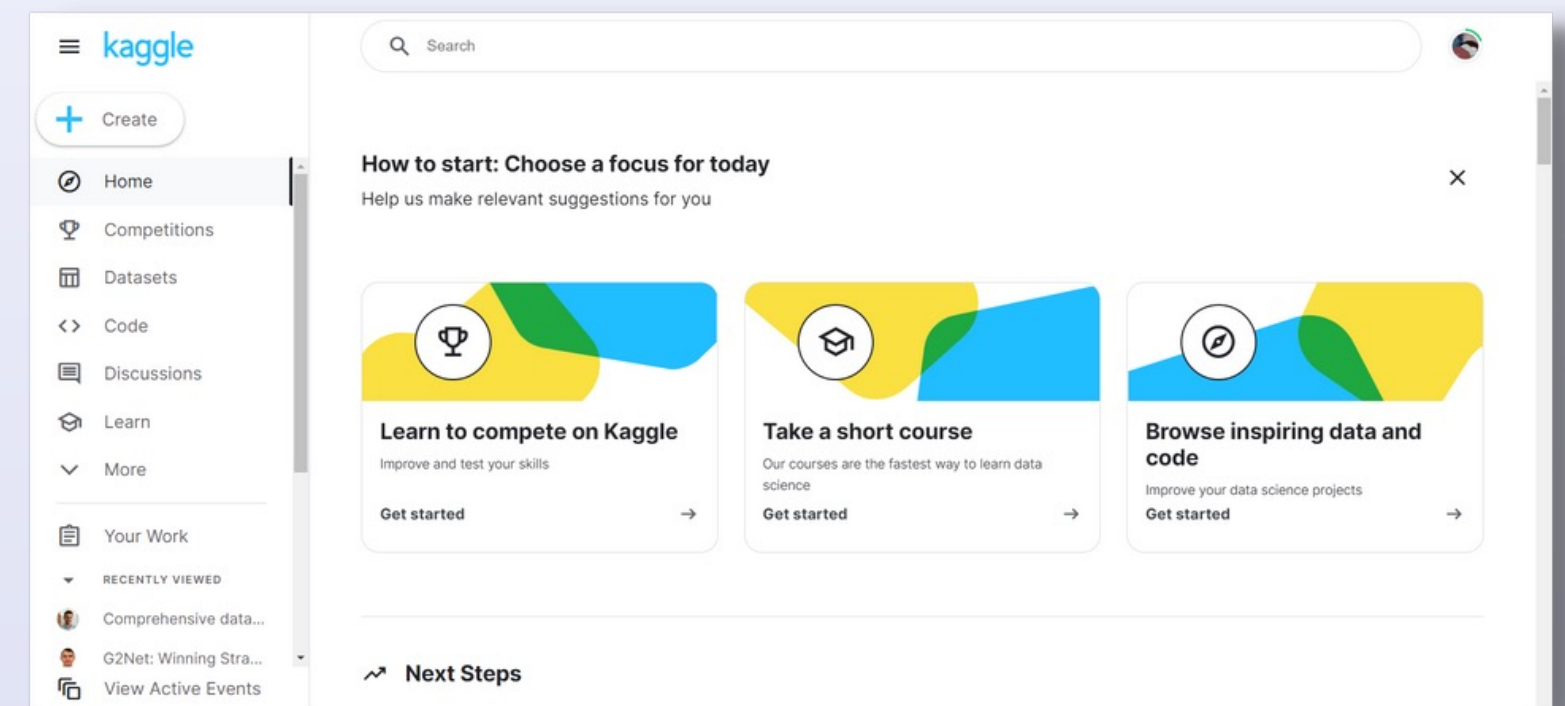
## FOR EXAMPLE:

**MACHINE LEARNING**

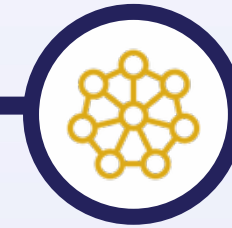**COMPETE IN REAL-WORLD PROBLEMS**

**FIND OR PUBLISH DATASET**

**DISCUSS WITH OTHERS**

**Kaggle Web Page**

# Kaggle's User Tier

**Novice**

A new user who joins Kaggle

**Contributor**

A user who has completed a profile engaged with the community and fully explored the platform of Kaggle.

**Expert**

A user who receives 5 bronze medals

**Master**

A user who receives 10 silver medals

**Grandmaster**

A user who receives 15 gold medals

# CODE RECOMMENDATION

Code recommendation helped improve developer productivity significantly.

Such tools have been gaining a lot of attention.
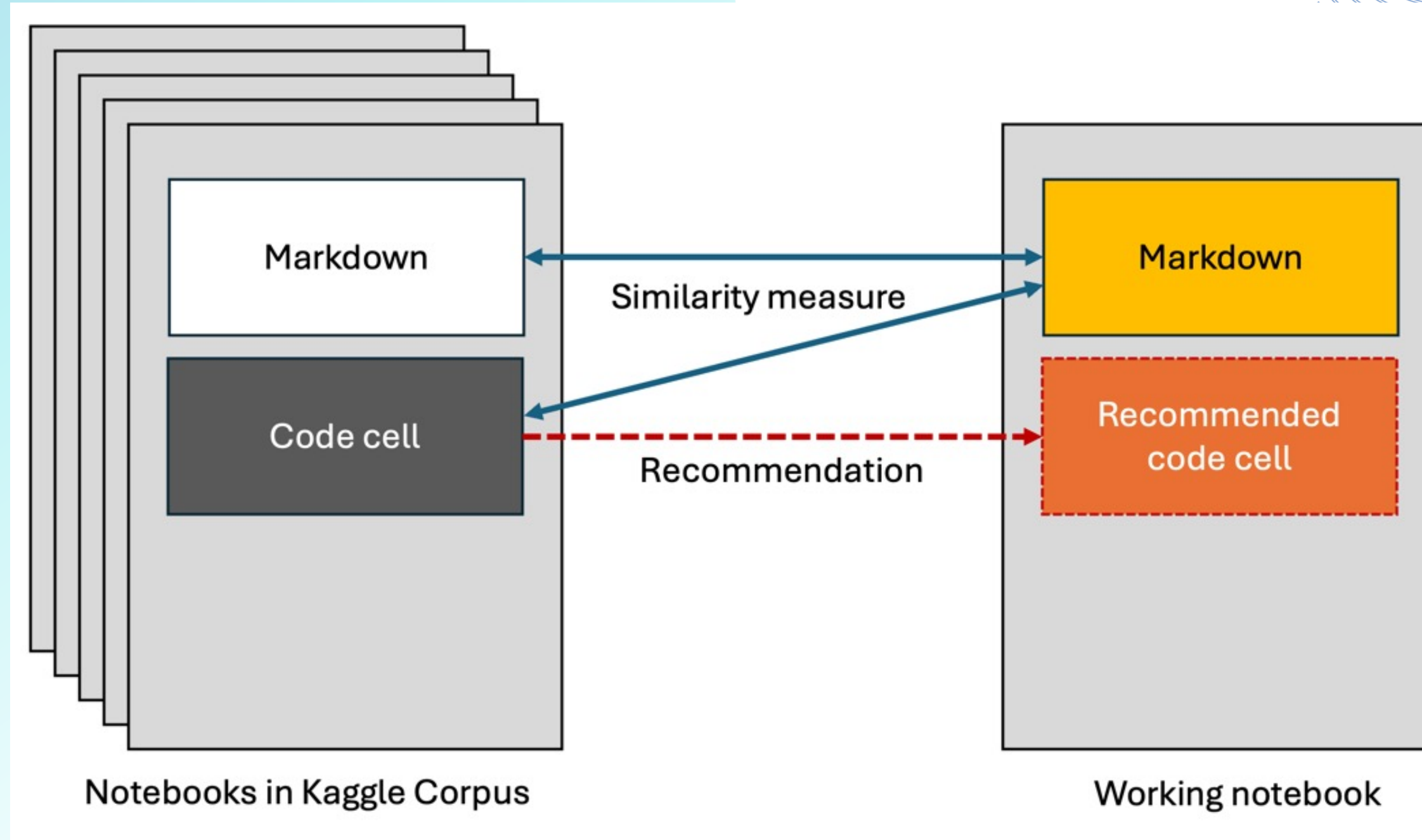
Many new tools are AI-powered

# EXISTING CODE RECOMMENDATION TECHNIQUES

| Name | Input | Approach |
| --- | --- | --- |
| Example Overflow | Text | Similarity of keywords search on database using fine-tuned TF-IDF weight |
| Copilot | Text/Code | OpenAI Codex model trained on large open-source projects in GitHub |
| Tabnine | Text/Code | AI-based proprietary algorithm |
| Aroma | Code | Similarity distance from code using clustering and ranking code snippets |
| Strathcona | Code | Similarity from user's local structural detail and code structural detail in repositories |
| Senatus | Code | Similarity of query input and indexed code using Minhash-LSH technique |

Can we recommend a code cell based on the given markdown cell by searching from existing Jupyter notebooks?
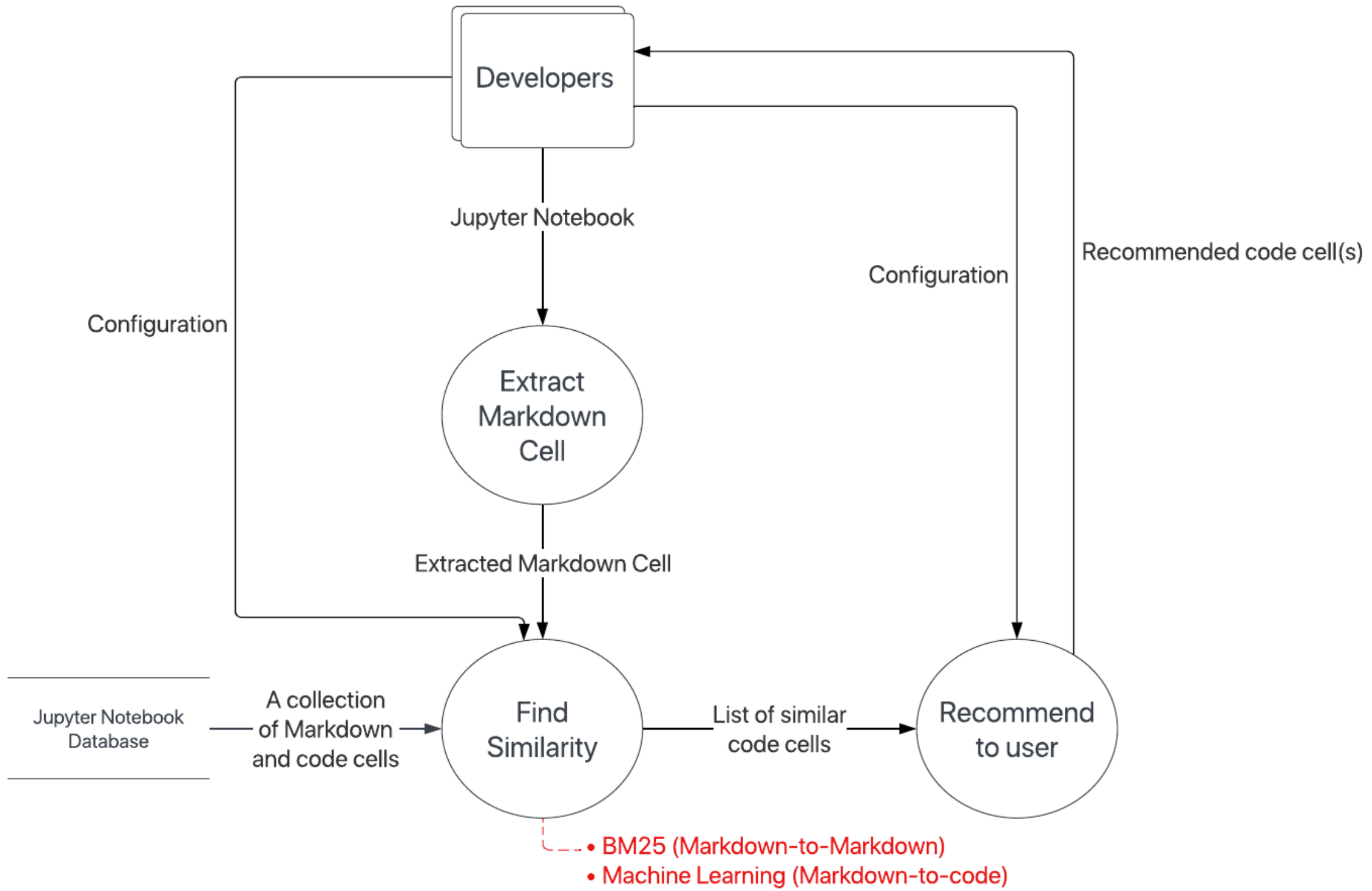
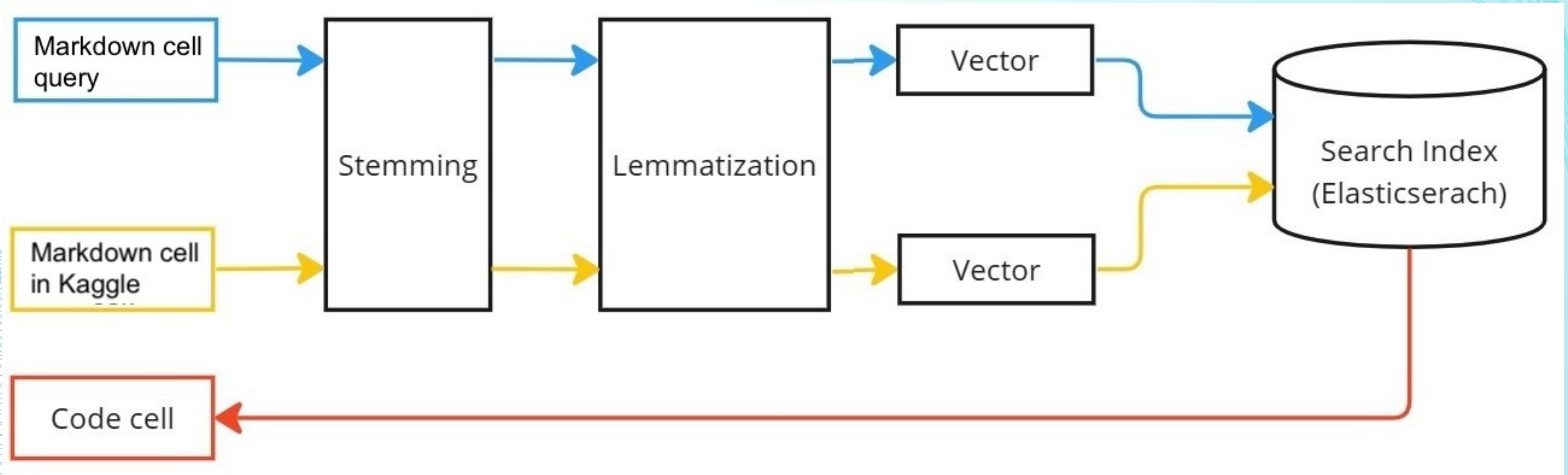# OUR APPROACH

# TYPHON ANALYSIS PROCESS

# SIMILARITY TECHNIQUES

# BM25

BM25 algorithm is a bag-of-words retrieval function, which ranks a set of documents regarding the query terms appearing in each document regardless of the proximity within the document.

BM25 has been widely used in search engines, such as Elasticsearch, as it is a robust and effective way to rank documents by relevance.

$$\sum_{i}^{n} IDF(q_i) \frac{f(q_i, D) \ * \ (k1 + 1)}{f(q_i, D) \ + \ k1 \ * \ (1 \ - \ b \ + \ b \ * \ \frac{fieldLen}{avgFieldLen})}$$
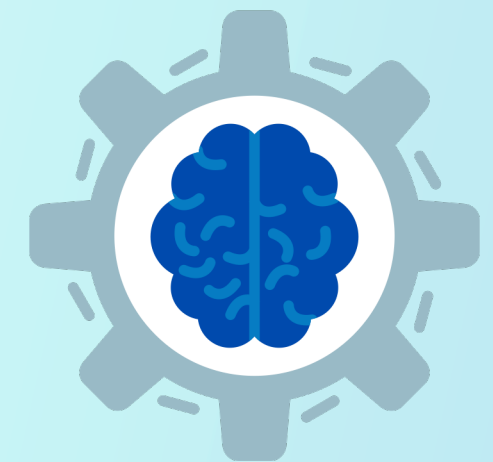
# Typhon with BM25
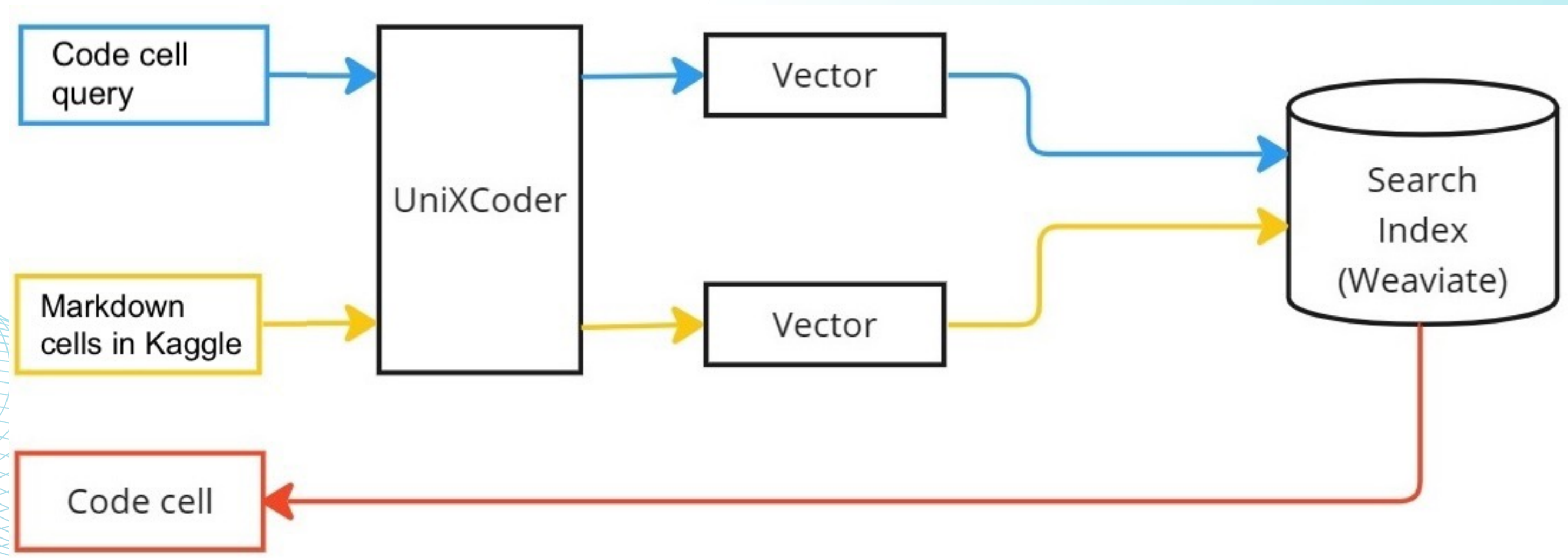


Text < -- > Text ----→ Code

# UniXcoder

A unified **cross-modal pre-trained model for the programming language**, which specializes in code understanding and code generation tasks.

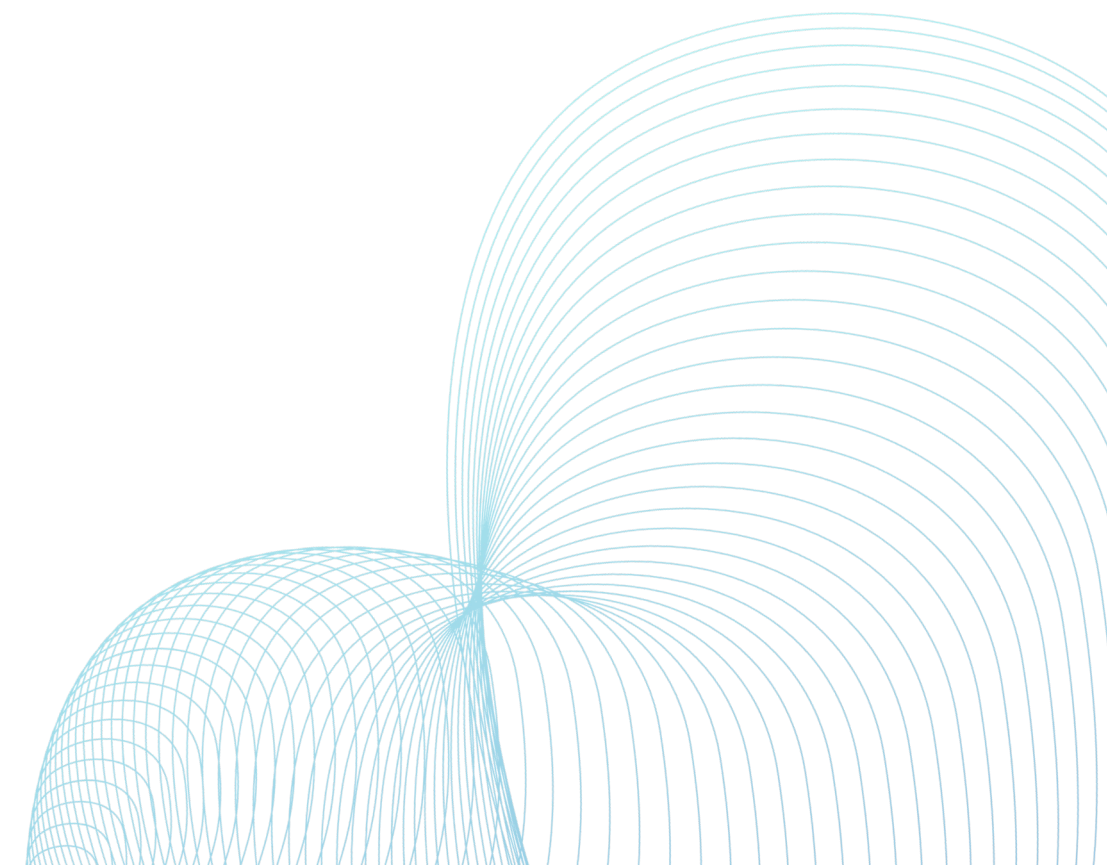It facilitates the usage of natural language and programming languages for code-related tasks.

Guo et al. (2022). UniXcoder: Unified Cross-Modal Pre-training for Code Representation. Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), 1, 7212–7225.

# Typhon with UniXcoder

# EVALUATION

# SANITY CHECK

| Rank | Type | Total Items | Total Correct | Total Correct (%) |
|---|---|---|---|---|
| Grand Master | UniXCoder | | 254 | 10.01 |
| | BM25 | 2,517 | 2,399 | 95.31 |
| | BM25 + stemming and lemmatization | | 2,132 | 84.72 |
| Master | UniXCoder | | 377 | 10.07 |
| | BM25 | 3,744 | 3,391 | 90.57 |
| | BM25 + stemming and lemmatization | | 3,007 | 80.32 |
| Expert | UniXCoder | | 605 | 6.33 |
| | BM25 | 9,553 | 8,644 | 90.48 |
| | BM25 + stemming and lemmatization | | 7,193 | 75.30 |

# Matplotlib Visualization Code Cell Recommendation

| Plot type | Sub plot type | Query term |
|---|---|---|
| Basic | Scatter | plot data using scatter visualization |
| | Bar | plot data using bar visualization |
| | Stem | plot data using stem visualization |
| | Step | plot data using step visualization |
| | Fill_between | plot data using fill_between visualization |
| | Stackplot | plot data using stackplot visualization |
| Plots of Arrays and Fields | Imshow | plot data using imshow visualization |
| | Pcolormesh | plot data using pcolormesh visualization |
| | Contour | plot data using contour visualization |
| | Contourf | plot data using contourf visualization |
| | Barbs | plot data using barbs visualization |
| | Quiver | plot data using quiver visualization |
| | Streamplot | plot data using streamplot visualization |
| Statistics Plots | Hist | plot data using hist visualization |
| | Boxplot | plot data using boxplot visualization |
| | Errorbar | plot data using errorbar visualization |
| | Violinplot | plot data using violinplot visualization |
| | Eventplot | plot data using eventplot visualization |
| | Hist2d | plot data using hist2d visualization |
| | Hexbin | plot data using hexbin visualization |
| | Pie | plot data using pie visualization |
| Unstructured Coordinates | Tricontour | plot data using tricontour visualization |
| | Tricontourf | plot data using tricontourf visualization |
| | Tripcolor | plot data using tripcolor visualization |
| | Triplot | plot data using triplot visualization |
| 3D | 3D Scatterplot | plot data using 3D scatterplot visualization |
| | 3D Surface | plot data using 3D surface visualization |
| | Triangular 3D Surface | plot data using triangular 3D surface visualization |
| | 3D Voxel , Volumetric Plot | plot data using 3D voxel , volumetric plot visualization |
| | 3D Wireframe Plot | plot data using 3D wireframe plot visualization |

https://matplotlib.org/stable/plot_types/index

# RESULTS

| Plot Type | Grand Master | | Master | | Expert | |
|---|---|---|---|---|---|---|
| | UniXcoder | BM25 | UniXcoder | BM25 | UniXcoder | BM25 |
| scatter | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| bar | ✓ | | ✓ | ✓ | ✓ | ✓ |
| step | ✓ | | | | ✓ | ✓ |
| imshow | ✓ | | ✓ | | | ✓ |
| contour | | | | | ✓ | |
| hist | | | | | | ✓ |
| boxplot | | ✓ | | | ✓ | ✓ |
| errorbar | | | ✓ | | | |
| violinplot | | | | ✓ | ✓ | |
| pie | ✓ | | ✓ | ✓ | | |
| tripcolor | ✓ | | | | | |
| 3d scatterplot | | ✓ | ✓ | | ✓ | ✓ |
| 3d surface | | | ✓ | | ✓ | |
| triangular 3d surface | | | | | ✓ | |
| Total Correct | 6 | 3 | 7 | 4 | 9 | 7 |
| Precision | 0.43 | 0.21 | 0.50 | 28.57 | 64.29 | 0.50 |

# TYPHON VS CODE EXTENSION

# CONCLUSION Ty

We propose **Typhon**

- An approach for recommending code cells based on existing code cells from Jupyter notebooks in the Kaggle dataset.

- We investigate using BM25 and UniXcoder for code and text similarity measurement.

- We performed an evaluation based on matplotlib visualizations and found moderate accuracy in recommendations with UniXcoder outperforming BM25.

- Our Typhon VS code extension is available in the marketplace.