# Md. Muhimenul Mubin
# Student ID : 20101112
# Cse431 Task 2
# Section : 01

Paper name : ***ERNIE-Code: Beyond English-Centric Cross-lingual Pretraining for Programming Languages***

## Introduction to the Paper

ERNIE-Code is a language model that connects programming languages and natural languages in a universal multilingual fashion. It goes beyond English-centric pre-training for programming languages and covers 116 natural languages and 6 programming languages. ERNIE-Code outperforms previous models in various code intelligence tasks, including code summarization, codegeneration, and code retrieval. The model achieves strong performance in both zero-shot and few-shot settings, and it can generate code in multiple programming languages from natural language descriptions. The paper provides a detailed analysis of the model's performance and limitations, as well as suggestions for future work.

## Motivation and purpose

The motivation of this paper is to address the limitations of existing language models that are English-centric and do not adequately support multilingual programming languages. The purpose of this paper is to introduce ERNIE-Code, a unified pre-trained language model that bridges the gap between multilingual natural languages and programming languages. The authors aim to demonstrate the effectiveness of ERNIE-Code in various code intelligence tasks, including code summarization, code generation, and code retrieval, and to show its advantage in zero-shot and few-shot settings. The paper also provides insights into the challenges and opportunities of multilingual programming language modeling and suggests future directions for research.

## Dataset

The paper uses several public datasets for evaluating the performance of ERNIE-Code in various code intelligence tasks. These datasets include CodeSearchNet, CodeXGLUE, and the Multilingual Code Mixed Corpus. The authors use the same train-test splits for all downstream

tasks to ensure fair comparison across models. The paper provides detailed statistics of the evaluation datasets, including the number of samples, the number of programming languages, and the distribution of code lengths. The authors also report the evaluation metrics used in each task, such as accuracy, F1 score, and BLEU score. Finally, the paper includes an analysis of the strengths and limitations of the evaluation datasets and suggests ways to improve them for future research.

## Methodology

The methodology of this paper involves the development and evaluation of ERNIE-Code, a unified pre-trained language model for multilingual natural languages and programming languages. The authors employ two methods for universal cross-lingual pre-training: span-corruption language modeling and pivot-based translation language modeling. They use a large-scale corpus of multilingual text and code to pre-train ERNIE-Code and fine-tune it on various downstream tasks, including code summarization, code generation, documentation translation, and code repair. The authors use several evaluation metrics to assess the performance of ERNIE-Code, including accuracy, F1 score, and BLEU score. They also conduct ablation studies to analyze the contribution of different pre-training methods to the model's performance.

## Results and Analysis

The results and analysis of this paper show that ERNIE-Code outperforms previous multilingual language models for programming languages and natural languages across a wide range of code intelligence tasks. The authors demonstrate the effectiveness of ERNIE-Code in code summarization, code generation, documentation translation, and code repair tasks, and show its advantage in zero-shot and few-shot settings. The paper provides detailed analysis of the model's performance on each task, including the evaluation metrics used and the comparison with baseline models. The authors also conduct ablation studies to analyze the contribution of different pre-training methods to the model's performance. The paper includes qualitative findings that demonstrate the model's ability to bridge the semantics and syntax between multilingual natural language instructions and programming language functions. Finally, the authors provide insights into the strengths and limitations of ERNIE-Code and suggest future directions for research.

## Limitations and Future work:
So far, the primary constraint I've identified is that this paper focuses on brief text documents and the dataset is relatively small compared to other NLP classification studies. By using superior datasets it will be easier to get superior results.