# Assignment II: Machine Translation

Weijian Deng ( kendwj@hku.hk )

Shu Chen ( schen59@hku.hk )

Find it on Github: https://github.com/hkukend/DASC7606A-B-A2

# Optional Assignment:

- **This assignment is optional** due to the coming final exams

- **The final assignment grades** will be the maximum between Assignment 1 and 2

# Objectives:

- **Use the Hugging Face Hub for datasets:** discover, download, and prepare translation datasets
- **Load pretrained models and tokenizers** from the community
- **Build training pipelines using** transformers, datasets, accelerate, and evaluate
- **Fine-tune a translation model** end-to-end for a chosen language pair (zh-sim->en)
- **Evaluate translation quality** using BLEU (via sacrebleu) and report results
- **Develop debugging skills** for identifying and fixing common deep learning issues in HF version

# HuggingFace🤗:

# HuggingFace🤗 [Datasets](:):



Use one line to load datasets from HugginFace platform:

```
>>> from datasets import load_dataset
>>> dataset = load_dataset("wmt19", "zh-en")
>>> dataset
DatasetDict({
    train: Dataset({
        features: ['translation'],
        num_rows: 25984574
    })
    validation: Dataset({
        features: ['translation'],
        num_rows: 3981
    })
})
```

# HuggingFace🤗 [Transformers](): 



Use one line to load tokenizers and pretrained models from HugginFace platform:
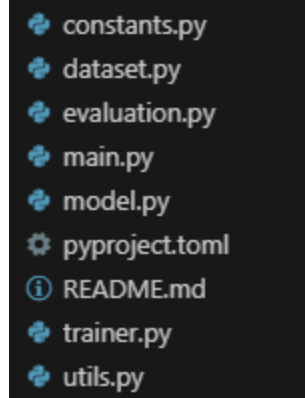
# What we will do:

## The workflow is organized into 5 files:

- **Dataset Sourcing & Preparation (dataset.py)**
  Explore HF Datasets, load splits, perform filtering/mapping, and train/validation/test preparation

- **Model & Tokenizer Setup (model.py)**
  Select a suitable base model from the Hub and initialize tokenizer and config

- **Training Pipeline (trainer.py)**
  Configure TrainingArguments, data collators, metrics, logging, and checkpointing

- **Run & Monitor Training (main.py)**
  Orchestrate end-to-end training and validation, with periodic evaluation

- **Evaluation & Reporting (evaluation.py)**
  Compute BLEU on the held-out test set; save artifacts and summary

constants.py
dataset.py
evaluation.py
main.py
model.py
pyproject.toml
README.md
trainer.py
utils.py

# What we will do:

- **Goal**: Fine-tune a suitable HF model on a translation dataset and achieve competitive BLEU.
- **What you can modify**: You may change any files in the repo **except** `main.py`, part of `evaluation.py` and `utils.py`. The test set must not be modified.
- **What to improve**:
  - **Enriching Datasets**: Use a more varied dataset (no leakage of the test set into training)
  - **Base model choice**: Select an appropriate pretrained model for your language pair
  - **Training pipeline**: Tune hyperparameters (batch size, LR, epochs, schedulers, label smoothing, gradient accumulation)
  - **Data processing**: Tokenization lengths, filtering, cleaning, language codes, special tokens
  - **Advanced HF features**: Mixed precision (`fp16`/`bf16`), gradient checkpointing, LoRA/PEFT, better data collators, scheduler choices, early stopping

## Example Accepted Datasets/Models

- Datasets: `wmt14`, `wmt16`, `wmt19`, `opus100`, `ted_talks_iwslt`, etc. (via HF Datasets)
- Models: MarianMT (Helsinki-NLP/opus-mt-xx-yy), mT5, MBART-50, M2M100, NLLB-200 (ensure your pair is supported), and even **LLMs**, etc.

---

- constants.py
- dataset.py
- evaluation.py
- main.py
- model.py
- pyproject.toml
- README.md
- trainer.py
- utils.py

# Grading:

We will re-run `main.py` and evaluate on the fixed test set. BLEU (SacreBLEU) is the primary metric.

**Important Considerations:**

1. **Error-Free Execution**: Your code must run without errors (and avoid GPU OOM on the provided environment)
2. **Correct Data Usage**: Do not alter or leak the test set into training
3. **No Personal Pre-trained Model**: "Downloads last month" of loaded pre-trained model on HuggingFace should be **GREATER THAN 10**
4. **Reasonable Performance**: Achieve competitive BLEU given the chosen model and setup
5. **Runtime**: Complete in a reasonable time budget (≤ 12 hours with **ONE GPU on HKU GPU Farm**)

**BLEU-based Grading (dummy thresholds, subject to adjustment):**

- **BLEU ≥ 25**: 100%
- **BLEU ≥ 24**: 90%
- **BLEU ≥ 23**: 80%
- **BLEU ≥ 22**: 70%
- **BLEU ≥ 21**: 60%
- **BLEU ≥ 20**: 50%
- **BLEU < 20 / Fail to reproduce / Overtime**: 0%

# Important dates:

- Assignment II Release: Nov. 06 (Thursday)

- Submission Deadline: Nov. 30 (Sunday) (23:59 GMT+8)

# Late submission policy :

- All submissions later than the deadline will NOT be accepted

# Questions!

If any more questions, please contact kendwj@hku.hk
or schen59@hku.hk