# Resource-Efficient NLP

KHIPU 2025

# What is Low-resource NLP?

- **Low-resource Data:** Dataset sizes which are smaller than usually available to train a given model.
  - This is going to be simulated in this tutorial using subsamples from Spanish/Portuguese datasets.

- **Low-resource Compute:** Constrained computing resources than usually required by standard training/inference methods.
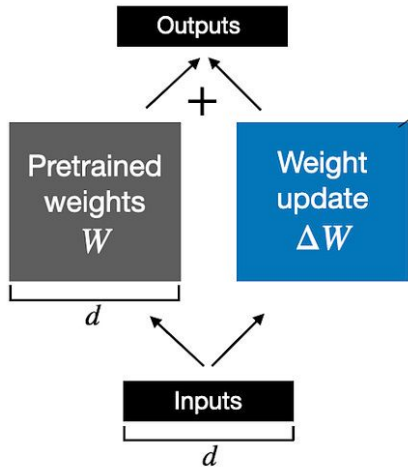
**Today:** Sampling of Resource-efficient NLP Toolkit

Today we'll cover a sample of different approaches which may be used to circumvent problems related to settings with low-resource data and/or compute!
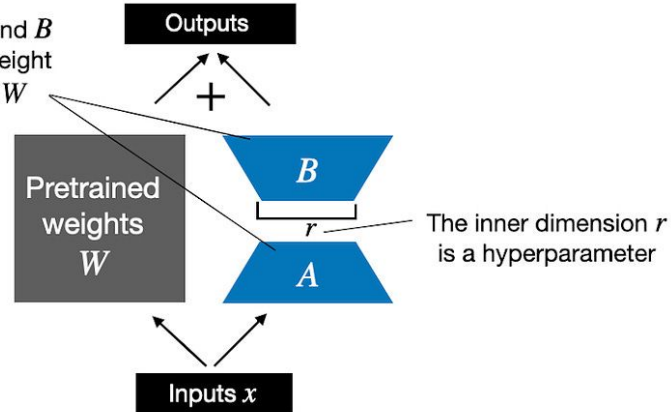
# Parameter-Efficient Fine-tuning (PEFT)

A family of methods which allow for modifying a subset of parameters during fine-tuning, reducing computational requirements.

# Pre-trained Models

Plenty of pre-trained models available open source in hubs like Hugging Face! This can save you compute resources that would otherwise be needed to pre-train a base model.
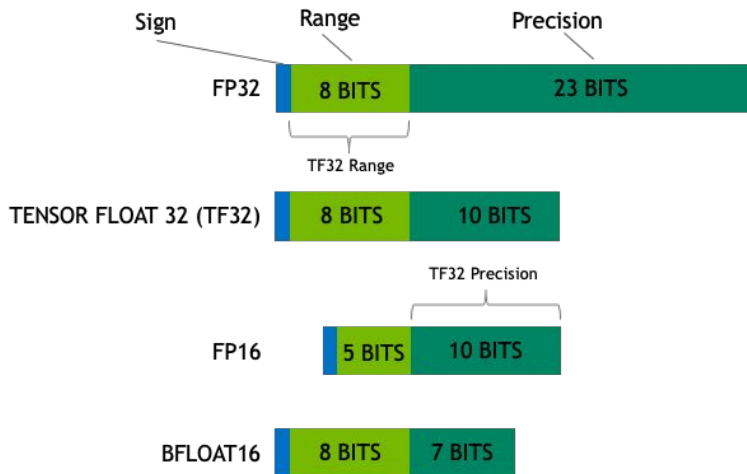
- Can use models off-the-shelf
- Can use models as a base to further fine-tune on your own data
- **Note:** Pre-trained models may not be readily available for low-resource data domains.

🤗 **huggingface_hub**

# Model Quantization

Model quantization methods reduce the precision of numerical representations of model weights and activations, yielding reductions in memory costs and potentially computational speedups as well.

# API Usage

Many state-of-the-art models are also available for use through APIs. Although APIs generally incur a cost to use, the cost is significantly lower than what would be required to train and house such models.