1
oneAPI
<HACK>ATHON

BUILD **SOLUTIONS** TO UNLOCK THE POTENTIAL OF **HETEROGENEOUS COMPUTING**

**LEAP** powered by Intel® oneAPI AI Analytics Toolkit

Problem Statement : **Open Innovation in Education**

Team Name : C5ailabs

Team Members : Rohit Sroch, Sujith R Kumar, Mohan K Rachumallu, Shubham Jain

# Problem Statement

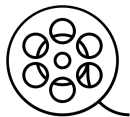| MOOCs (Massive Open Online Courses) | 200K Users in 2012 | 380M Users in 2020 | 34.26% CAGR 2022-27* | 5% -10% Completion rate |
|---|---|---|---|---|

## Key Challenges

**Lengthy videos**

**Instructor Availability**

**Slow response from forums**

**No real time Q&A/Mentor**

## Approach

### LEAP
(Learning Enhancement and Assistance Platform)

**AI based platform**

**Powered by Intel OneAPI**

intel. oneAPI BASE TOOLKIT

**Quality Education**

**All time Availability**

## Key Features of LEAP

**Ask Question/Doubt**

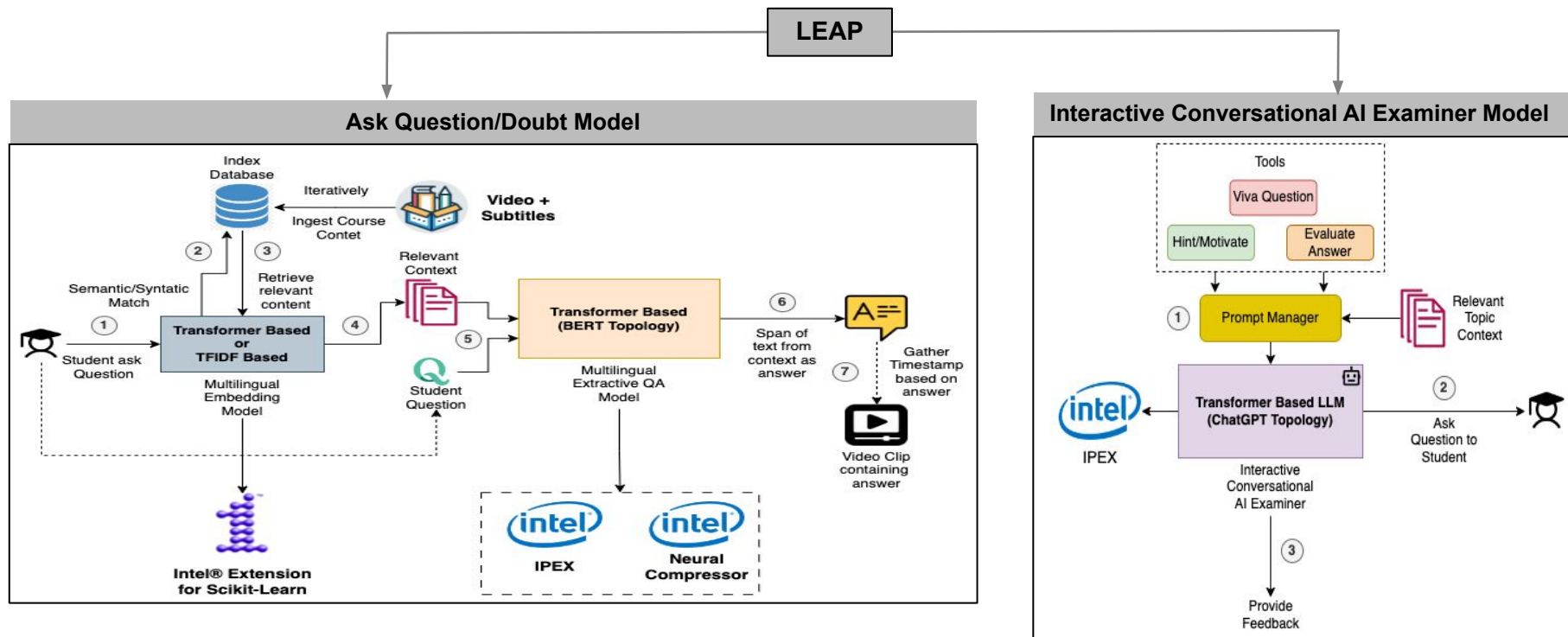**Conversational AI Examiner**

**Feedback from AI Examiner**

**Multilingual Support**

Reference: PRNewswire , Edtechreview; holonig

# High Level Architecture

Student can ask doubts

Student → Web UI → Select a course → Courses

**Courses**
- **Mathematics** → Algebra
- **Humanities** → Literature
- **Science** → Optics

Learning material

Video with Subtitle

Document

Ask Question 🎤 🔍

**NLP Models**
- Multilingual Embedding Model
- Multilingual Extractive QA Model

Predict → Video Clip containing Answer

Answer

Finished Learning

Interactive Conversational AI Examiner

AI Examiner conducts Viva exam of Student

Congraulations! on completing the Algebra topic. Here is 1st question for you

Solve, 4x + 5 when, x = 3

The answer is 17

Generative AI Model (ChatGPT like LLM)

Optimized Model using Intel® oneAPI

AI provides Student Feedback

Teacher

Optimized Models using Intel® oneAPI

**Intel® oneAPI AI Analytics Toolkit**

**Deep Learning**
- Intel® Extension For Pytorch
- Intel® Neural Compressor

**Machine Learning**
- Intel® Extension For Scikit-Learn

CPU | GPU | FPGA

**Here Intel® oneAPI AI Analytics Toolkit helped us to optimize the performance of LEAP by faster Inference and Training**

# LEAP: Detailed Model Architecture Diagram for Both Components

# Static-QAT-INT8 is **4.44x** times faster in performance as compared to baseline For our Extractive QA Model
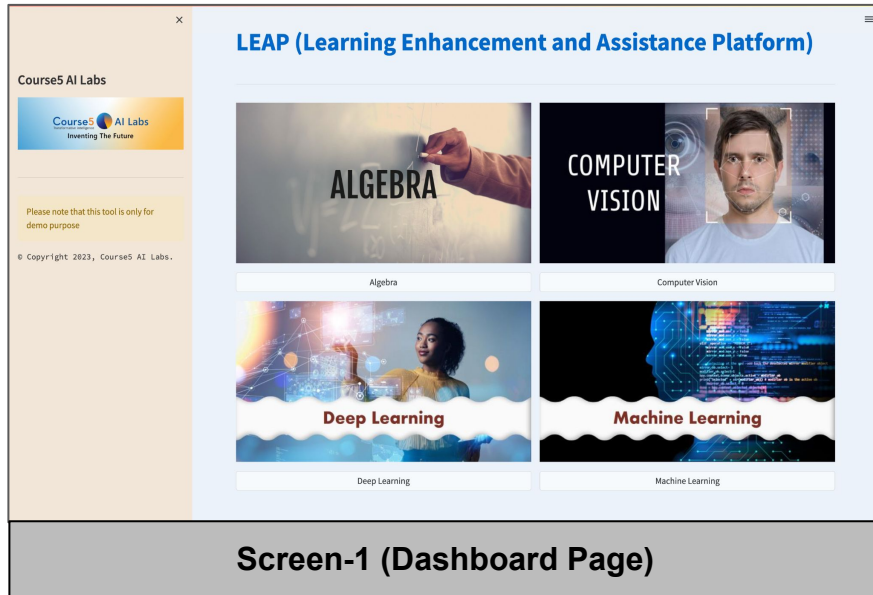
| Ask Question/Doubt Model | | | | | Interactive Conversational AI Examiner Model | | |
|---|---|---|---|---|---|---|---|

| Extractive Question Answering Model | | | | | TFIDF Embedding Model | | |
|---|---|---|---|---|---|---|---|
| | Pytorch (Base) - FP32 | Pytorch (IPEX) - FP32 | Static-QAT-INT8 | Static-Smooth-QAT-INT8 | | Scikit-Learn (Base) | Intel Extension For Scikit-Learn |
| Latency (milli sec) | 64.513 | 39.329 | **14.514** | 15.24 | Latency (milli sec) | 0.761 | **0.752** |
| Throughput (samples/sec) | 15.501 | 25.427 | **68.9** | 65.616 | Throughput (samples/sec) | 1313.63 | **1330.49** |
| F1 Score (SQuAD-v1) | **76.11** | **76.11** | 75.72 | 75.72 | | | |

*Table: Latency/Throughput/Speed-Up Benchmark result for **our Extractive Question Answering ALBERT Model (Multilingual) and TFIDF Embedding Model** on Intel® Dev Cloud machine (**Intel® Xeon® Platinum 8480+ (4th Gen: Sapphire Rapids) - 224v CPUs 503GB RAM**) with optimization using IPEX-FP32, Static-QAT-INT8 using Intel® Neural Compressor and TFIDFVectorizer using Intel® Extension for Scikit-Learn.*

# Demo Link and Screenshots

Link: https://www.youtube.com/watch?v=M51BFcoJa3k



Screen-1 (Dashboard Page)



Screen-2 (Week wise Content Page)

# Demo Screenshots



Screen-3 (Ask Question/Doubt Page)

Screen-4 (Ask Question/Doubt Page)

# Demo Screenshots



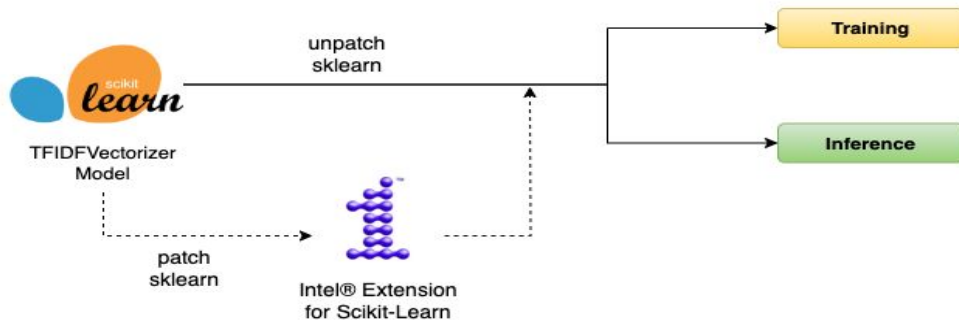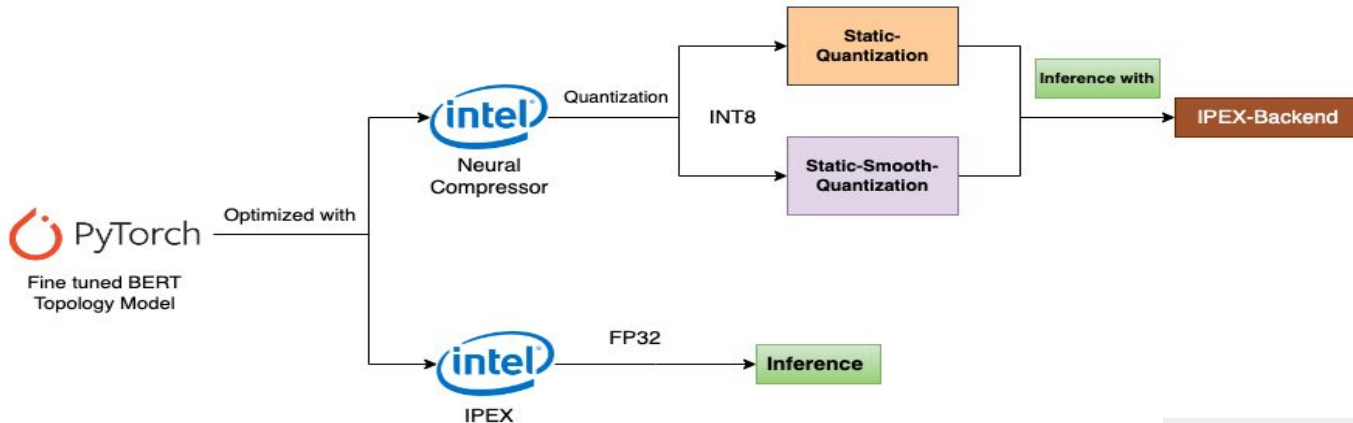Screen-5 (Interactive Conversational AI Examiner Asks Question to student)



Screen-6 (Interactive Conversational AI Examiner provides hints and motivates a student in case of a wrong answer)

**GitHub Link (Codes should be public and available after hackathon also)**

**https://github.com/rohitc5/intel-oneAPI**

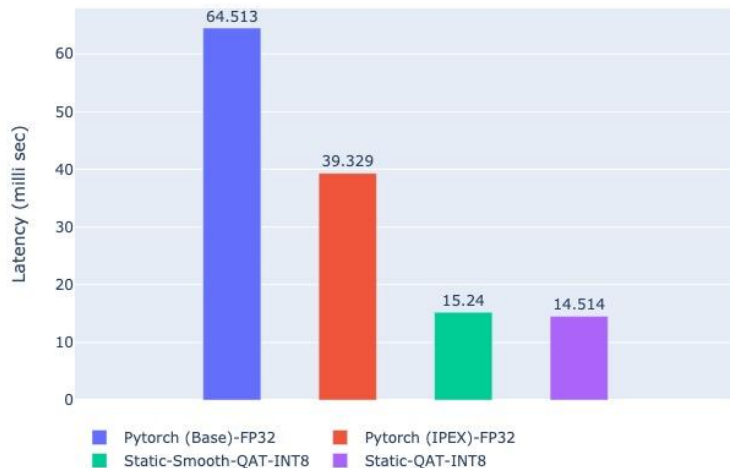# Result Summary (unique aspects of oneAPI/SYCL used)

Process flow to optimize the LEAP models by using

Fig 1. Intel® Extension for Pytorch, Intel® Neural Compressor and

Fig 2. Intel® Extension for Scikit-Learn

# Extractive QA Model (BERT Topology) Latency/Speed-Up Comparison with IPEX and Intel® Neural Compressor
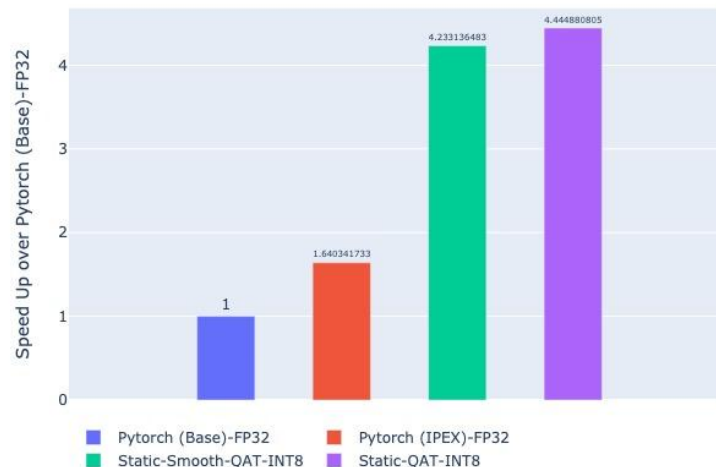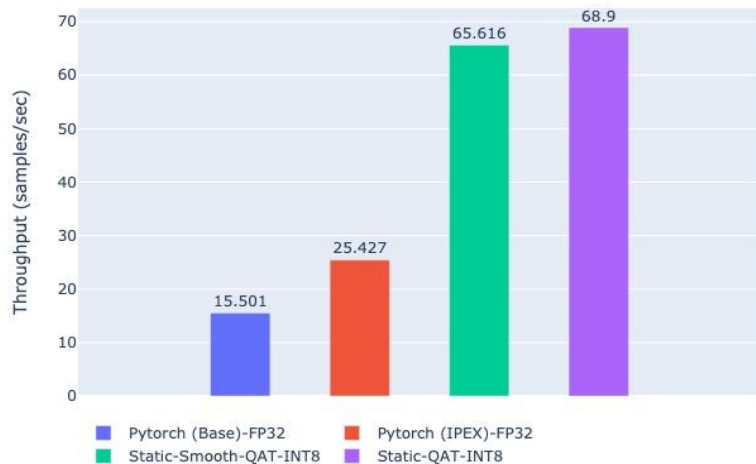


Fig: Latency/Speed-Up Benchmark result for **our Extractive Question Answering ALBERT Model (Multilingual)** on Intel® Dev Cloud machine (**Intel® Xeon® Platinum 8480+ (4th Gen: Sapphire Rapids) - 224v CPUs 503GB RAM**) with optimization using IPEX-FP32 and Static INT8-Quantization using Intel® Neural Compressor.

**For Ask Question/Doubt Extractive QA Model**

# Extractive QA Model (BERT Topology) Throughput/F1 Score Comparison with IPEX and Intel® Neural Compressor
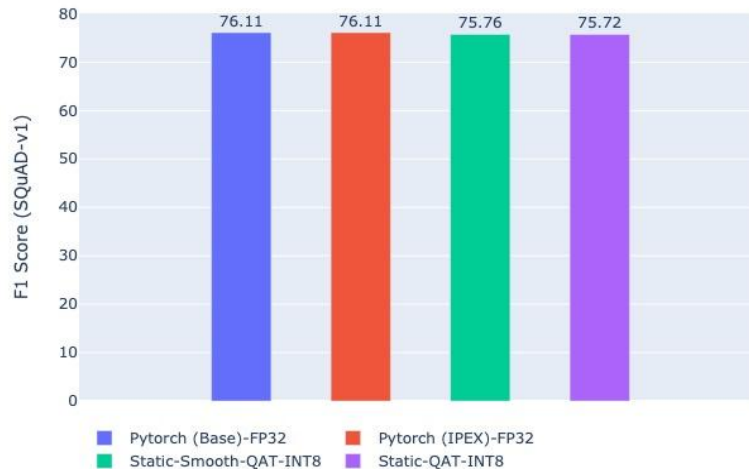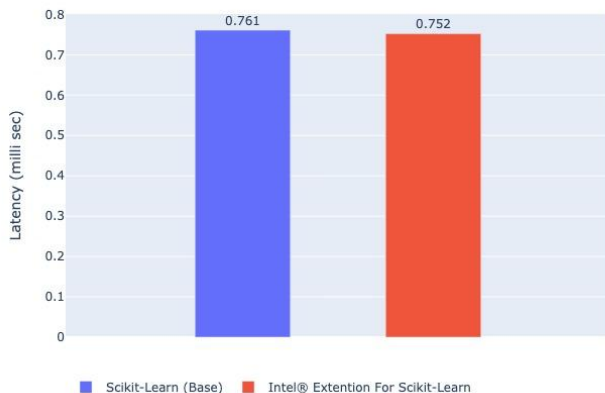
Fig: Throughput/F1 Score Benchmark result for **our Extractive Question Answering ALBERT Model (Multilingual)** on Intel® Dev Cloud machine (**Intel® Xeon® Platinum 8480+ (4th Gen: Sapphire Rapids) - 224v CPUs 503GB RAM**) with optimization using IPEX-FP32 and Static INT8-Quantiz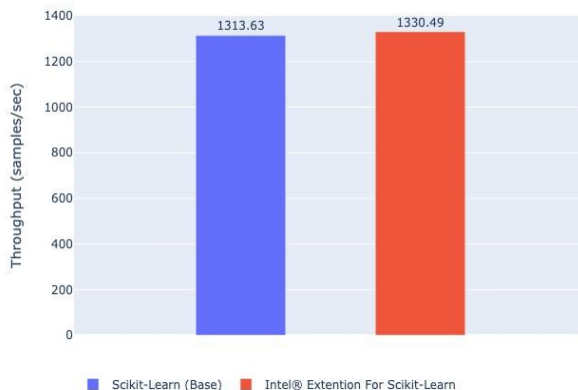ation using Intel® Neural Compressor. Also, the model (https://huggingface.co/ai4bharat/indic-bert) was fine-tuned on SQuAD-v1 dataset.

**For Ask Question/Doubt Extractive QA Model**
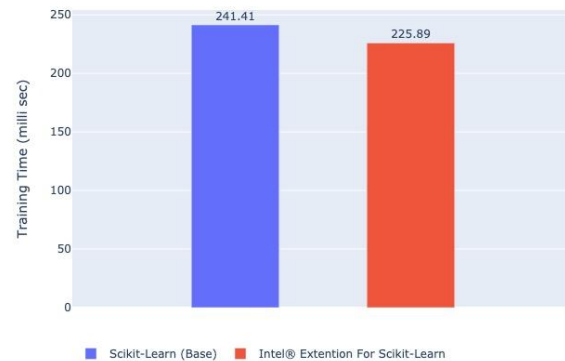
# Scikit-Learn (Base) vs Intel® Extension for Scikit-Learn



Fig: Benchmark results for **TFIDFVectorizer** Embedding model during training and inference on Intel® Dev Cloud machine (**Intel® Xeon® Platinum 8480+ (4th Gen: Sapphire Rapids) - 224v CPUs 503GB RAM**). Please Note that we don't see much of a difference may be because we used a tiny dataset.

For Ask Question/Doubt Embedding Model

# Model Checkpoint Release

https://huggingface.co/rohitsroch