**DAY2 REPORT SUMMARY ON TE:**

Transformer Engine (TE) is a performance-optimized library from NVIDIA that uses FP8 precision to speed up transformer-based deep learning models. TE is perfect for large-scale AI training and inference because it dynamically manages precision, which drastically lowers memory consumption and boosts throughput while preserving model correctness.

NVIDIA Hopper architecture GPUs, which have hardware capabilities especially built for effective low-precision processing, support full FP8 in Transformer Engine. This demonstrates how closely contemporary GPU hardware design and AI software optimisation are related.

An autocast context and an FP8 recipe are used in PyTorch to support FP8 computation. This enables the framework to automatically choose FP8 when it is advantageous and safe, dropping back to greater precision when necessary for numerical stability.

Transformer Engine is a high-performance, hardware-dependent library, making installation challenging. It frequently needs to be built against the system's GPU drivers and requires compatible NVIDIA GPUs as well as particular versions of CUDA, cuDNN, and PyTorch. While this guarantees optimal performance, installing TE requires more work than installing regular Python modules.