

Allocated GPU memory for regular and flash attention models

