

Average inference times for regular and flash attention models

