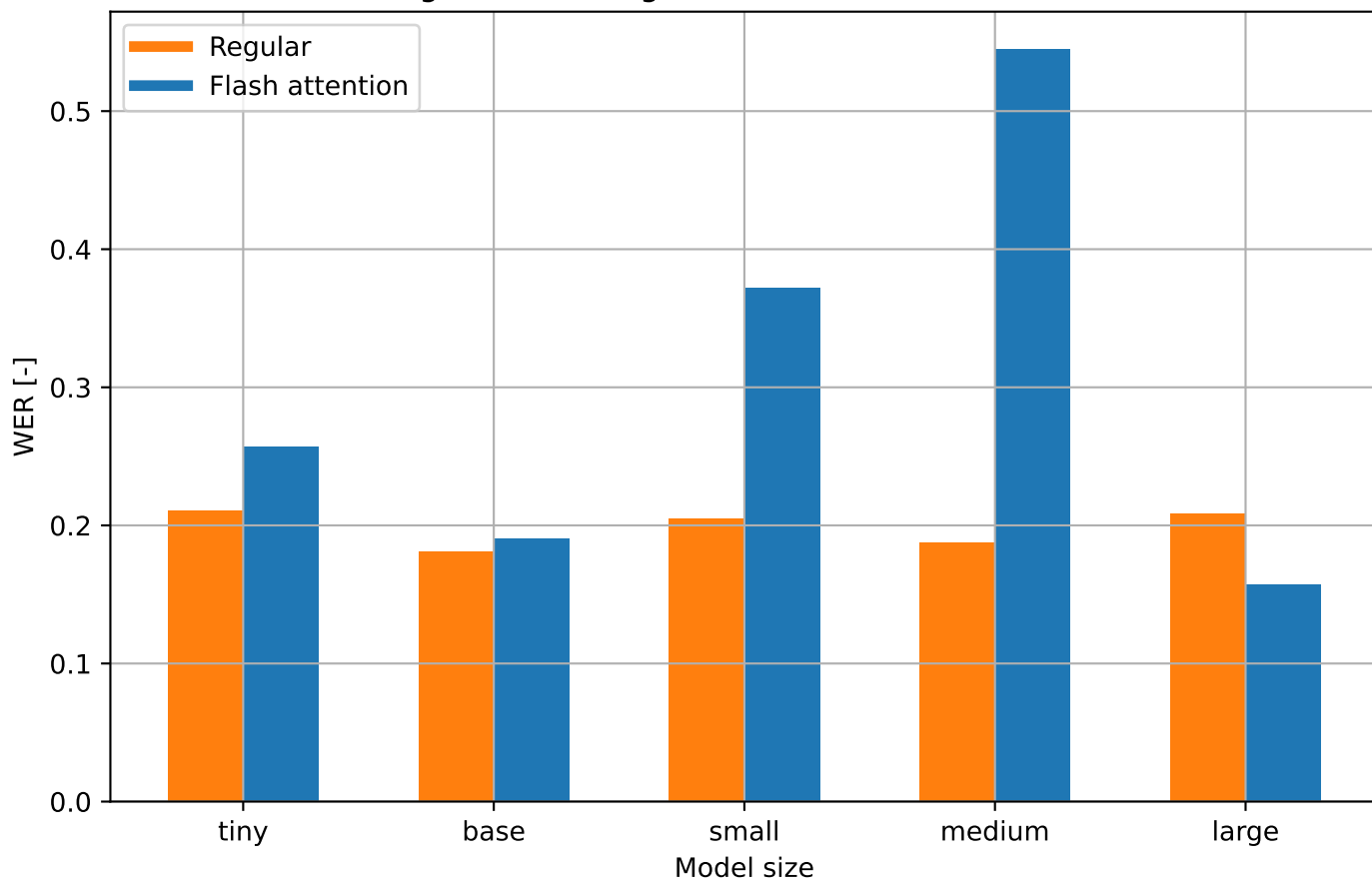
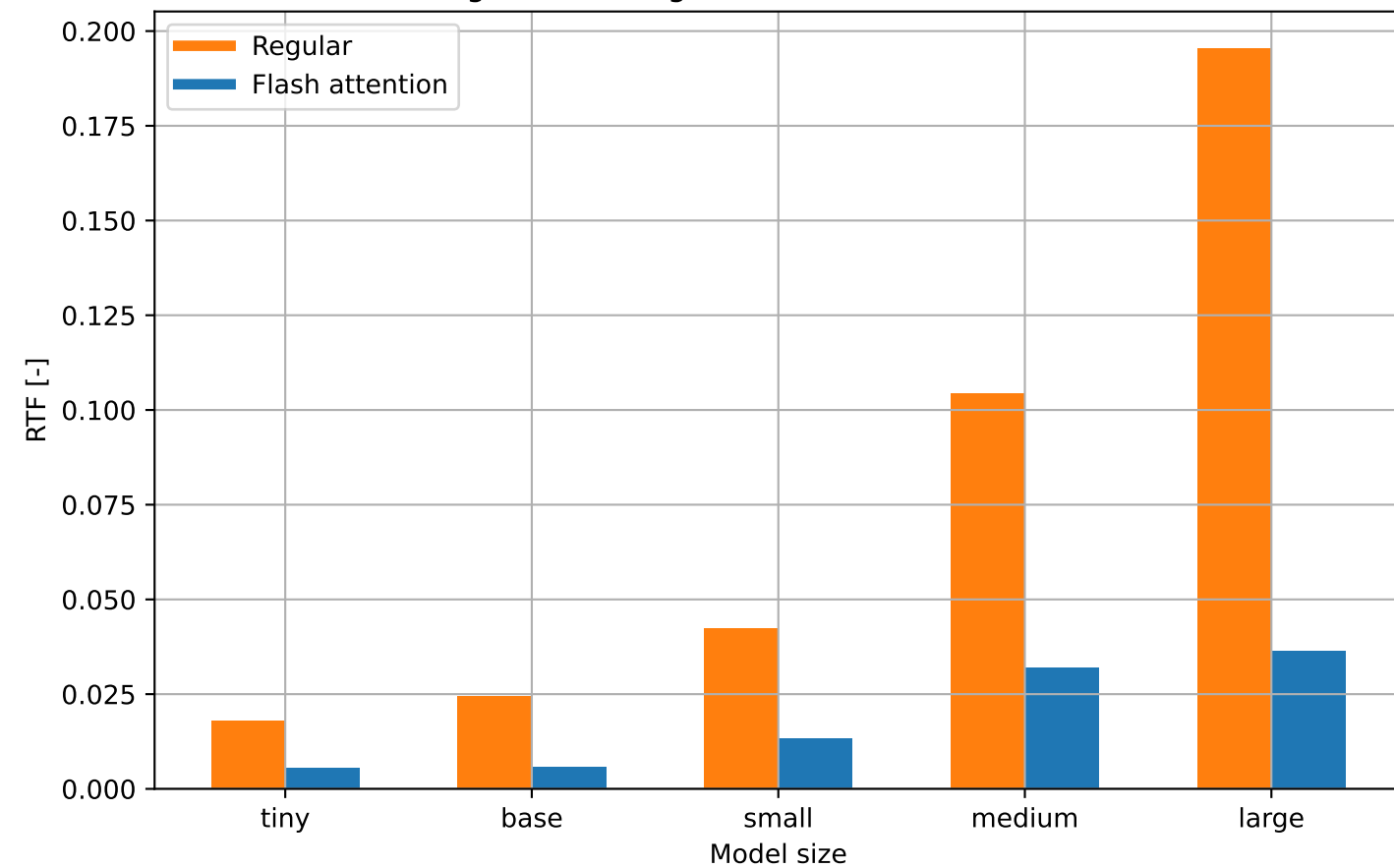


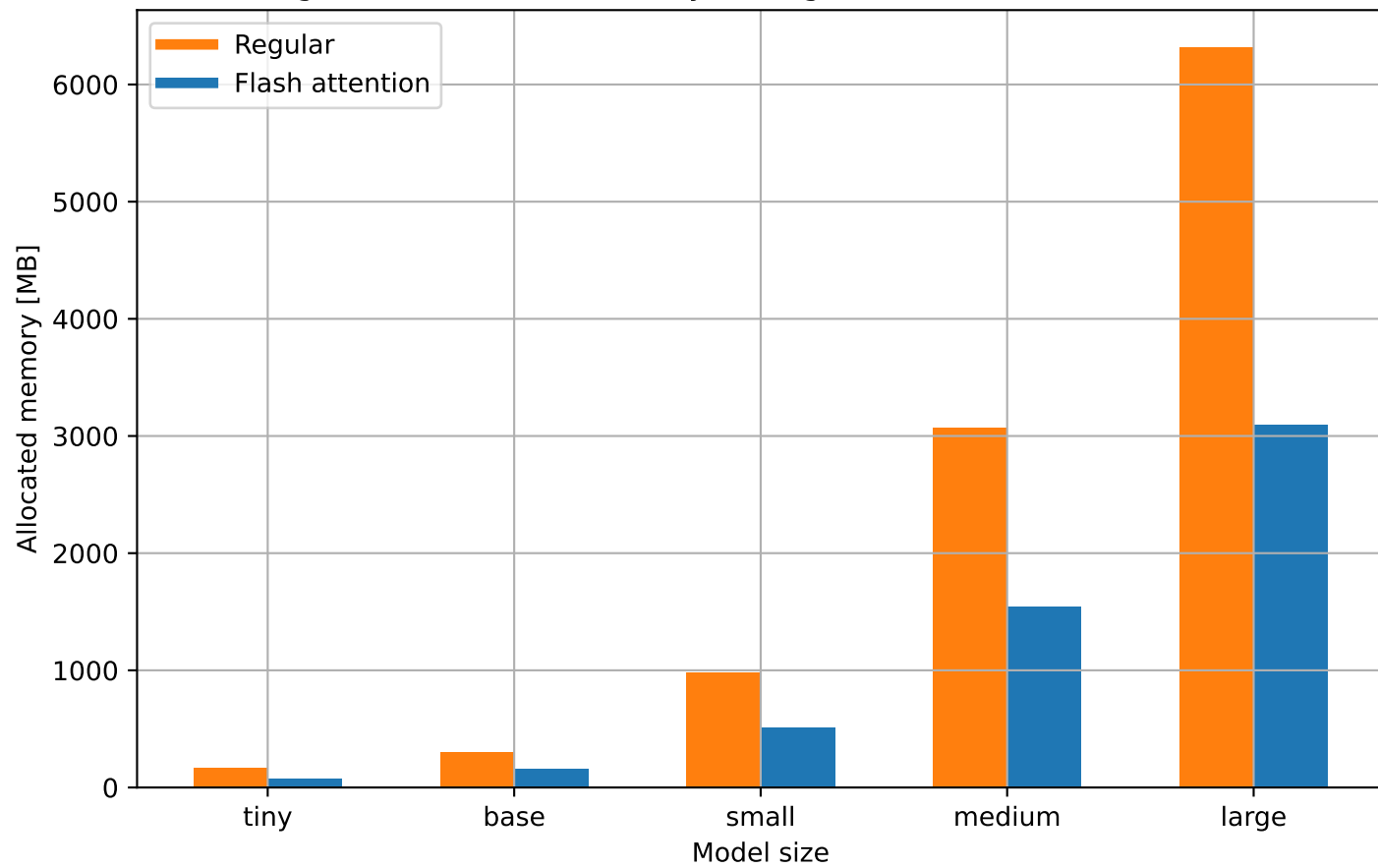
Average WER for regular and flash attention models



Average RTF for regular and flash attention models



Average allocated GPU memory for regular and flash attention models



Average inference time for regular and flash attention models

