

Average inference time for different deployment methods, model type: large

