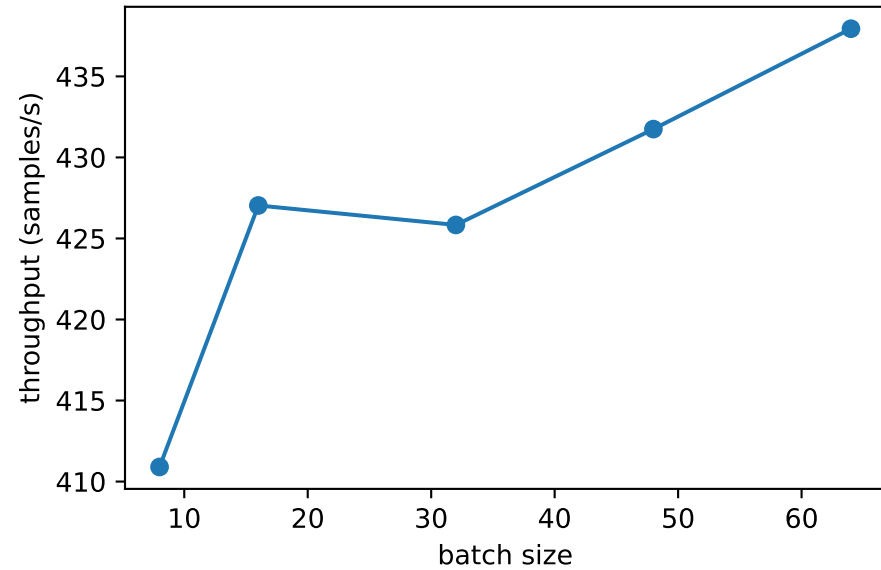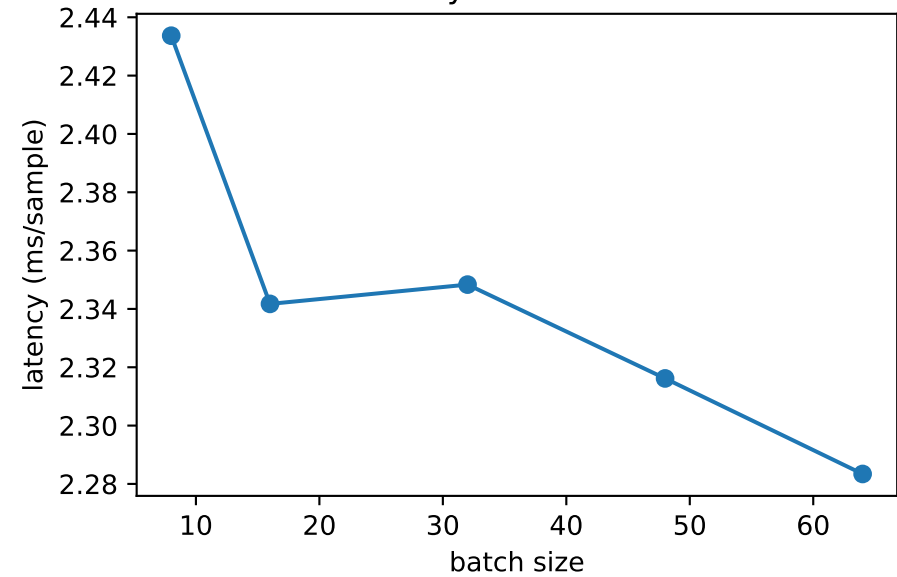Efficiency & resource curve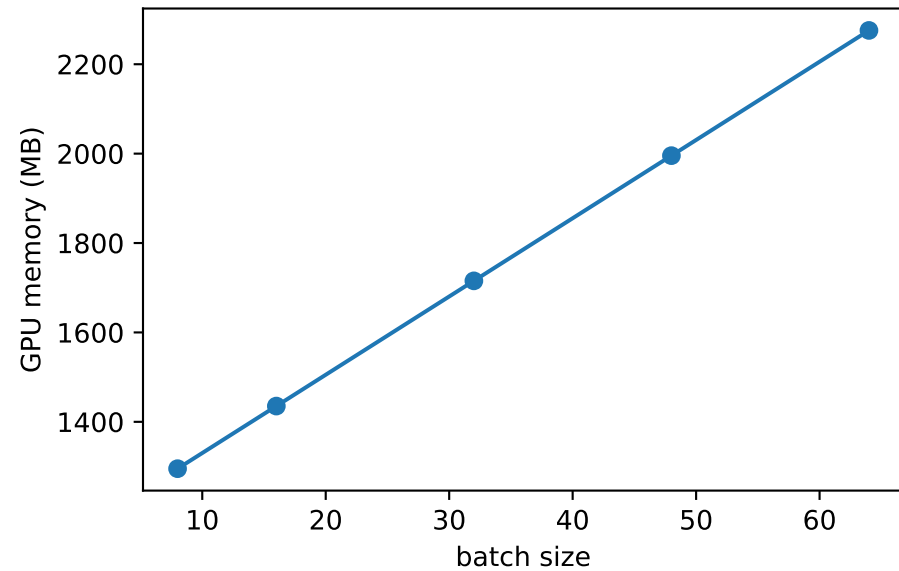s