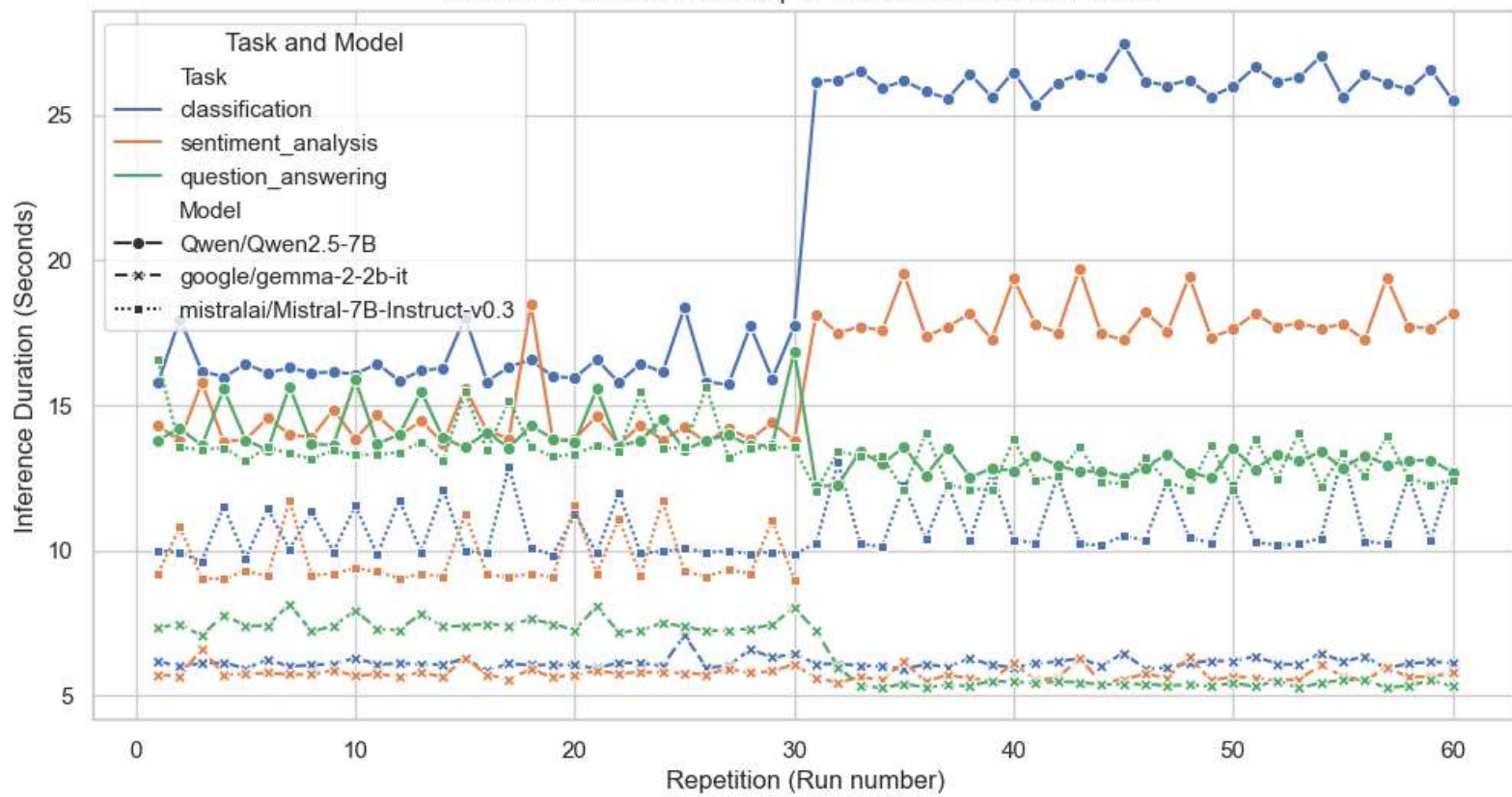
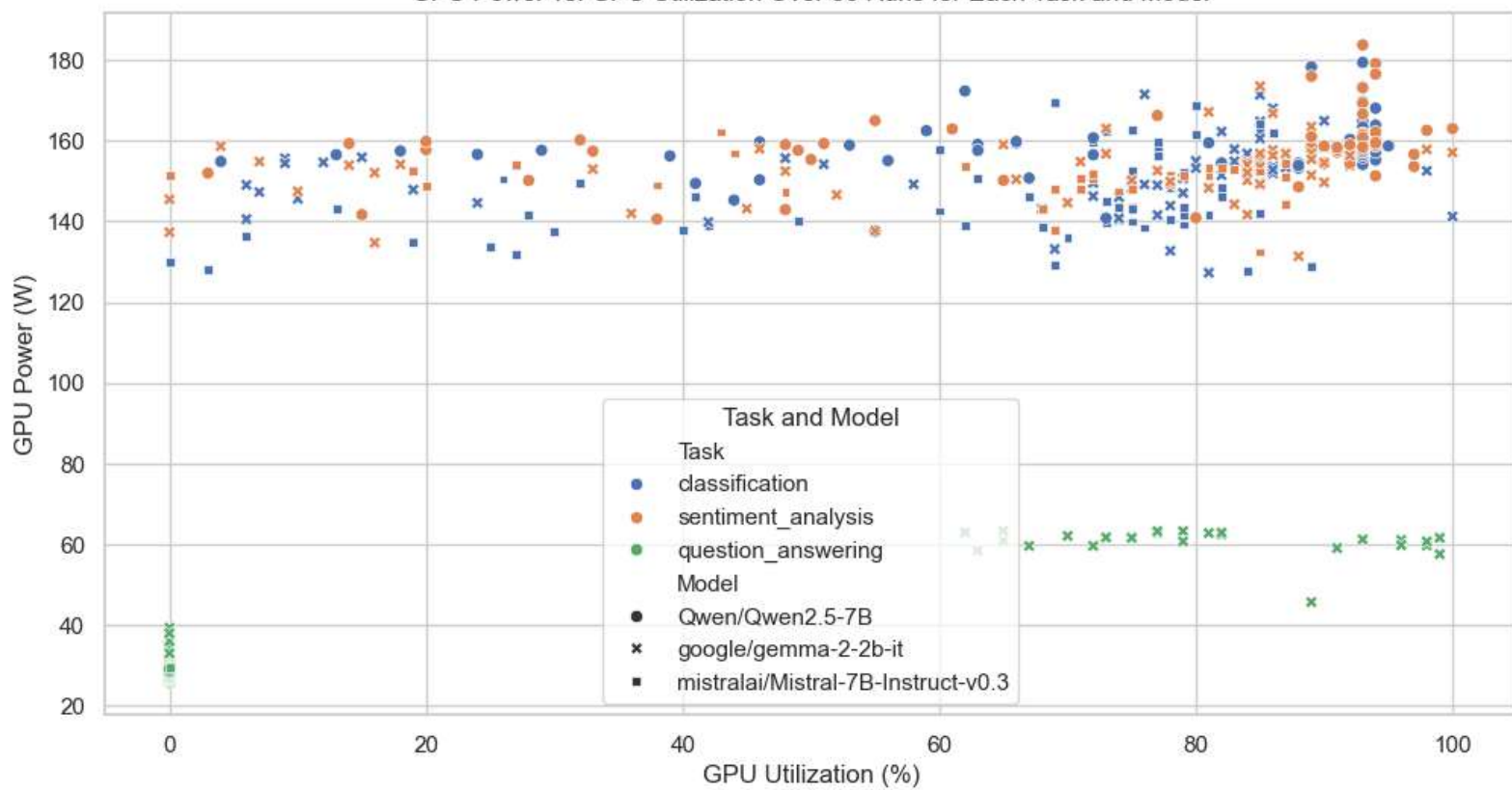


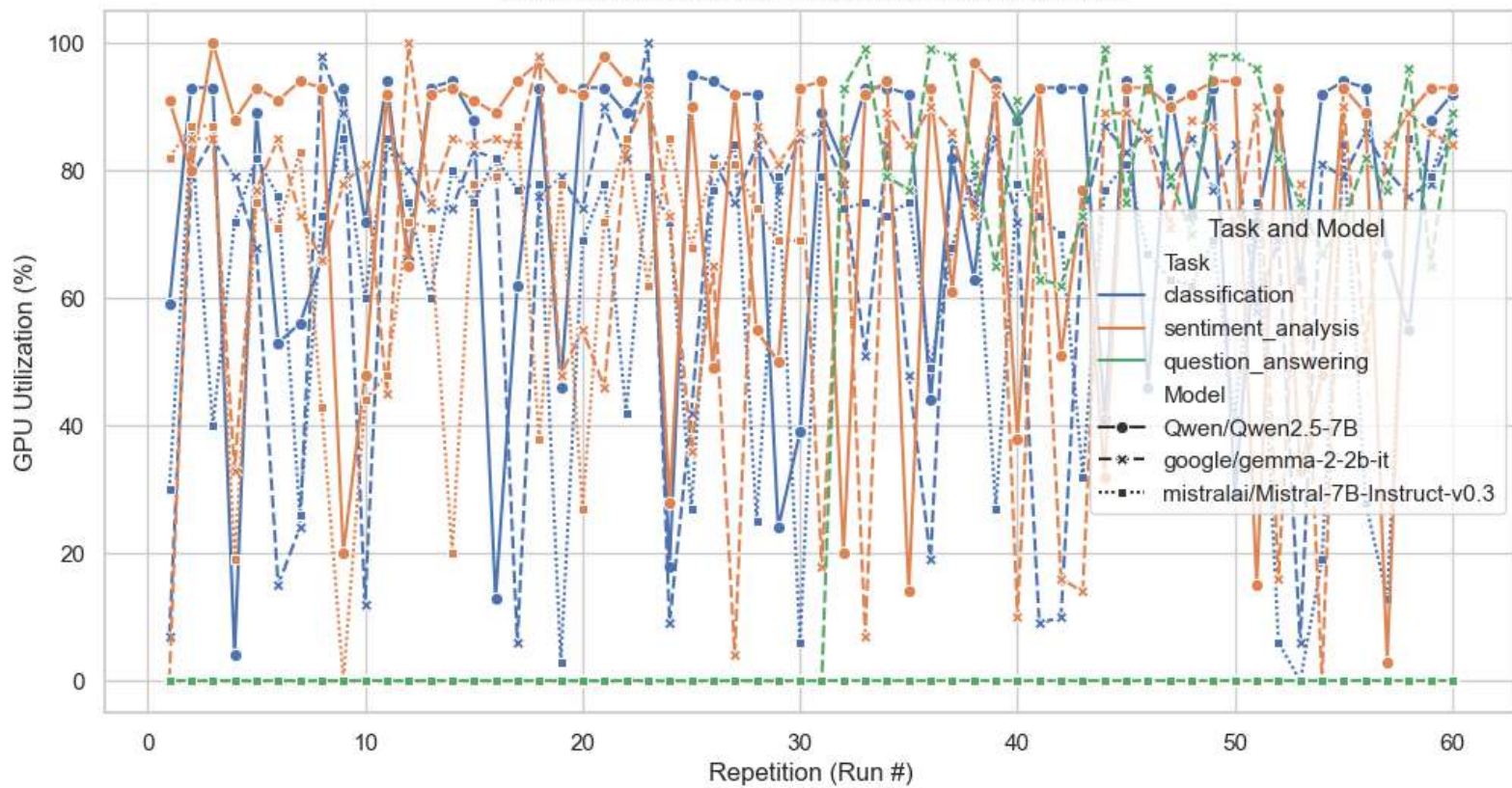
Inference Duration Across Repetitions for Each Task and Model

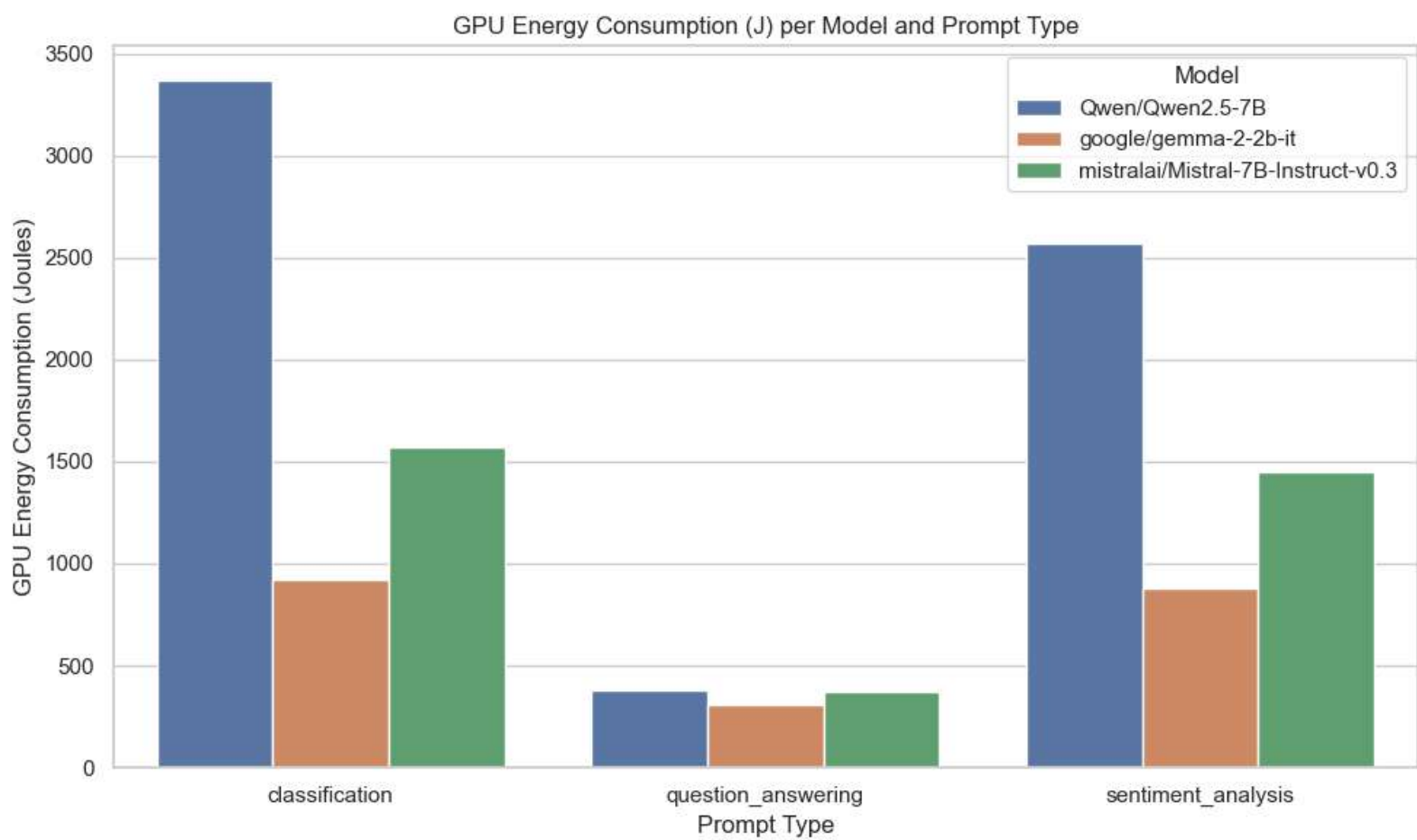


GPU Power vs. GPU Utilization Over 30 Runs for Each Task and Model

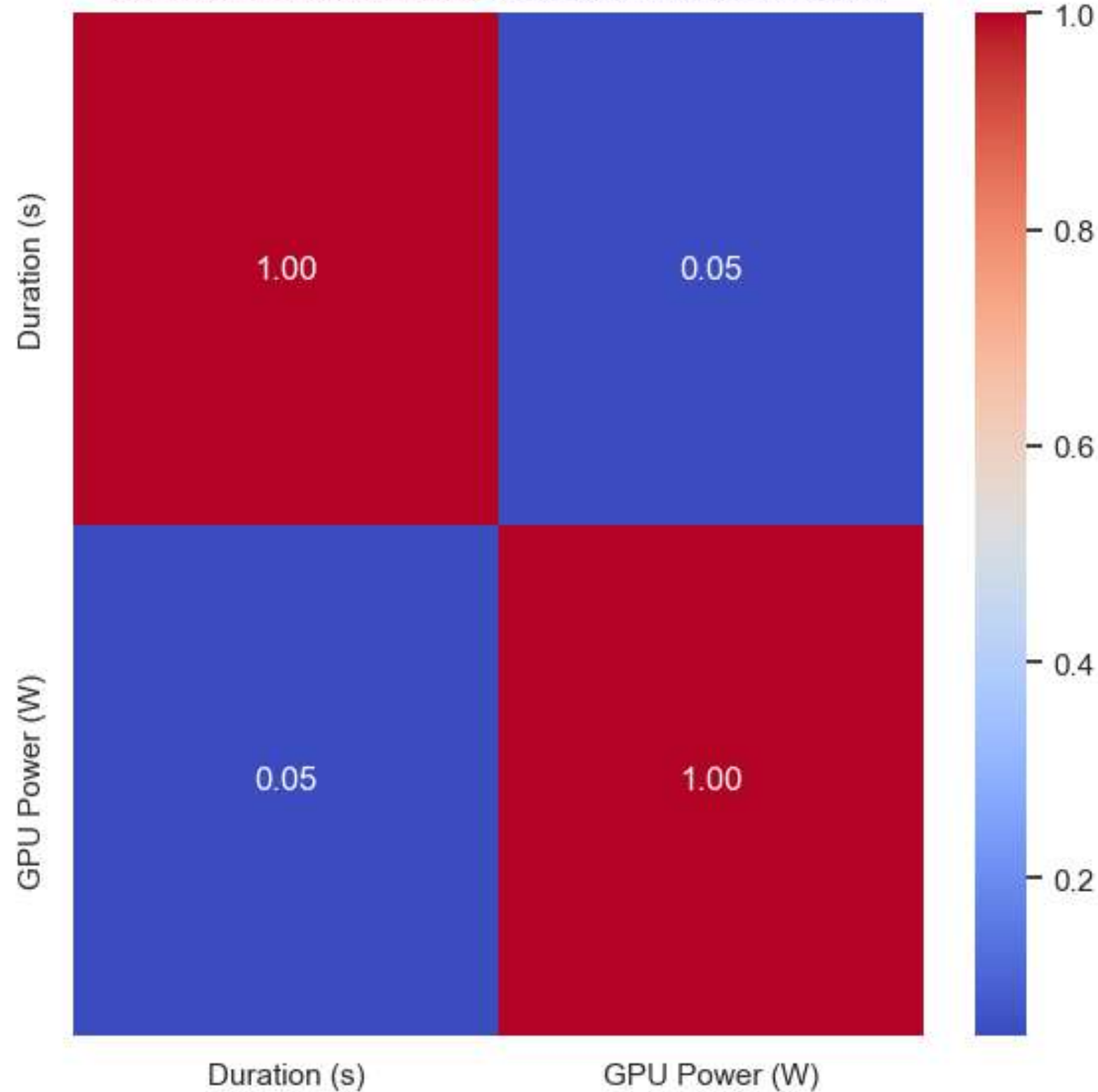


GPU Utilization Over 30 Runs for Each Task and Model

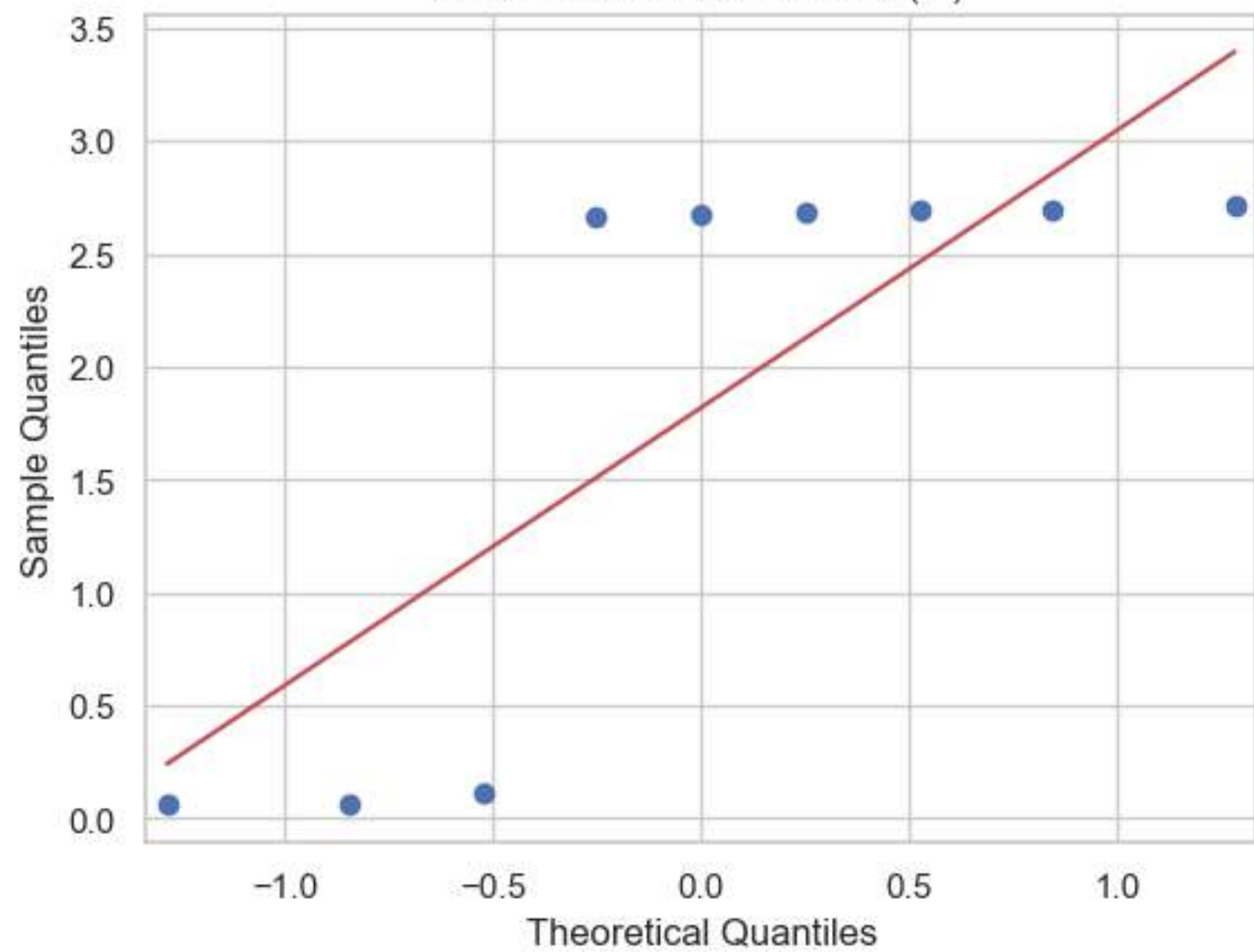




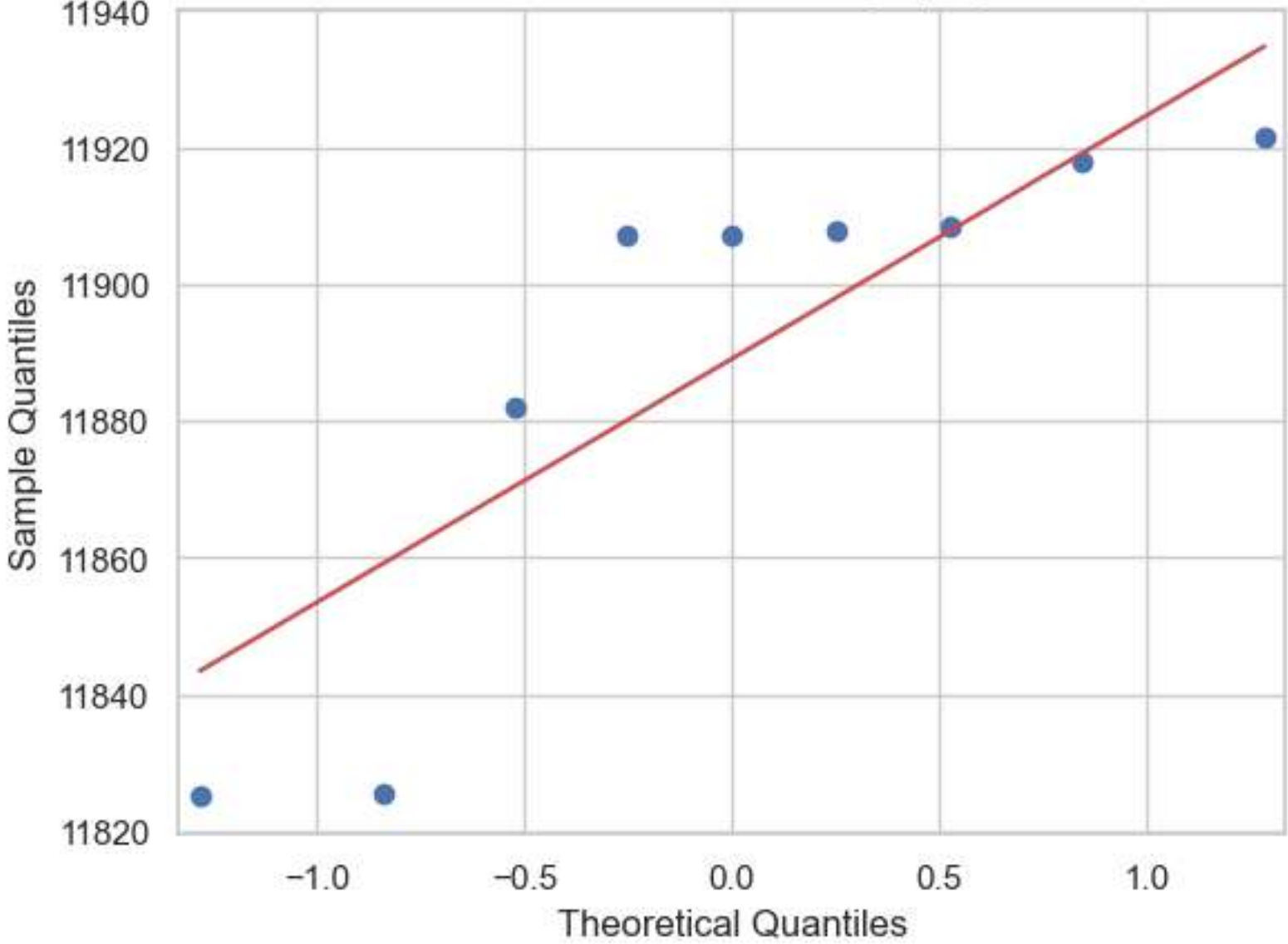
Correlation Between Inference Time and GPU Power



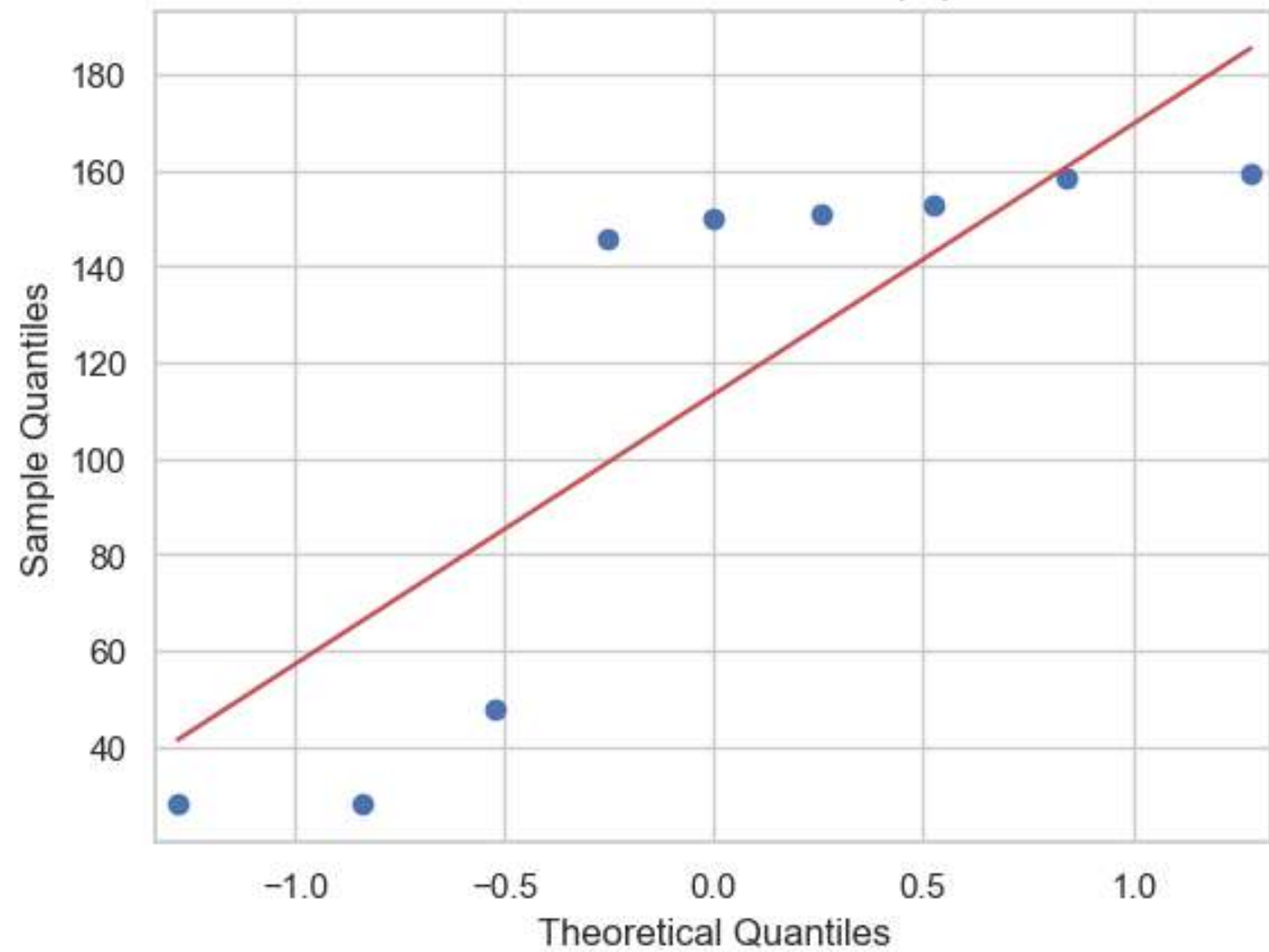
Q-Q Plot for CPU Utilization (%)



Q-Q Plot for VRAM Usage (MB)

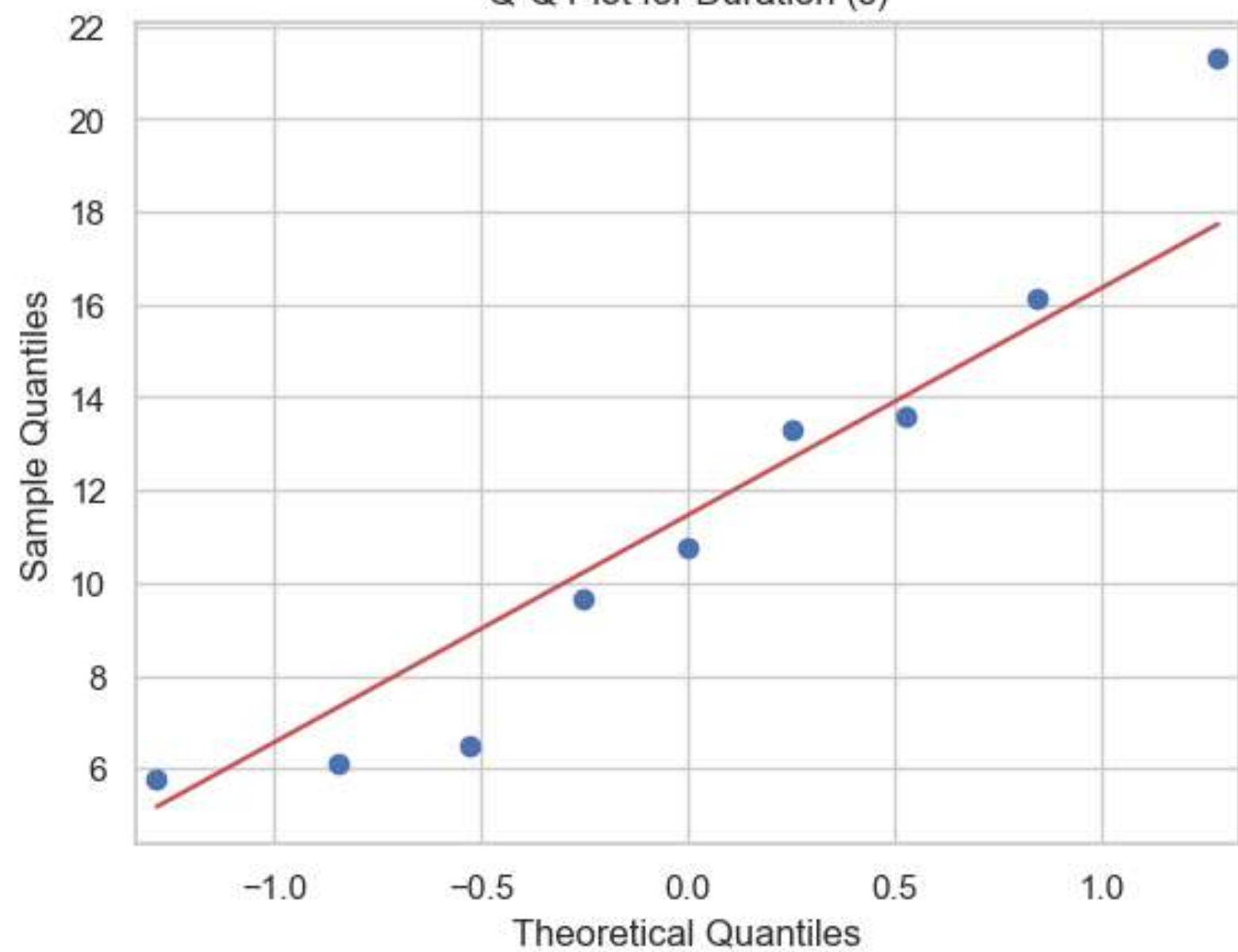


Q-Q Plot for GPU Power (W)

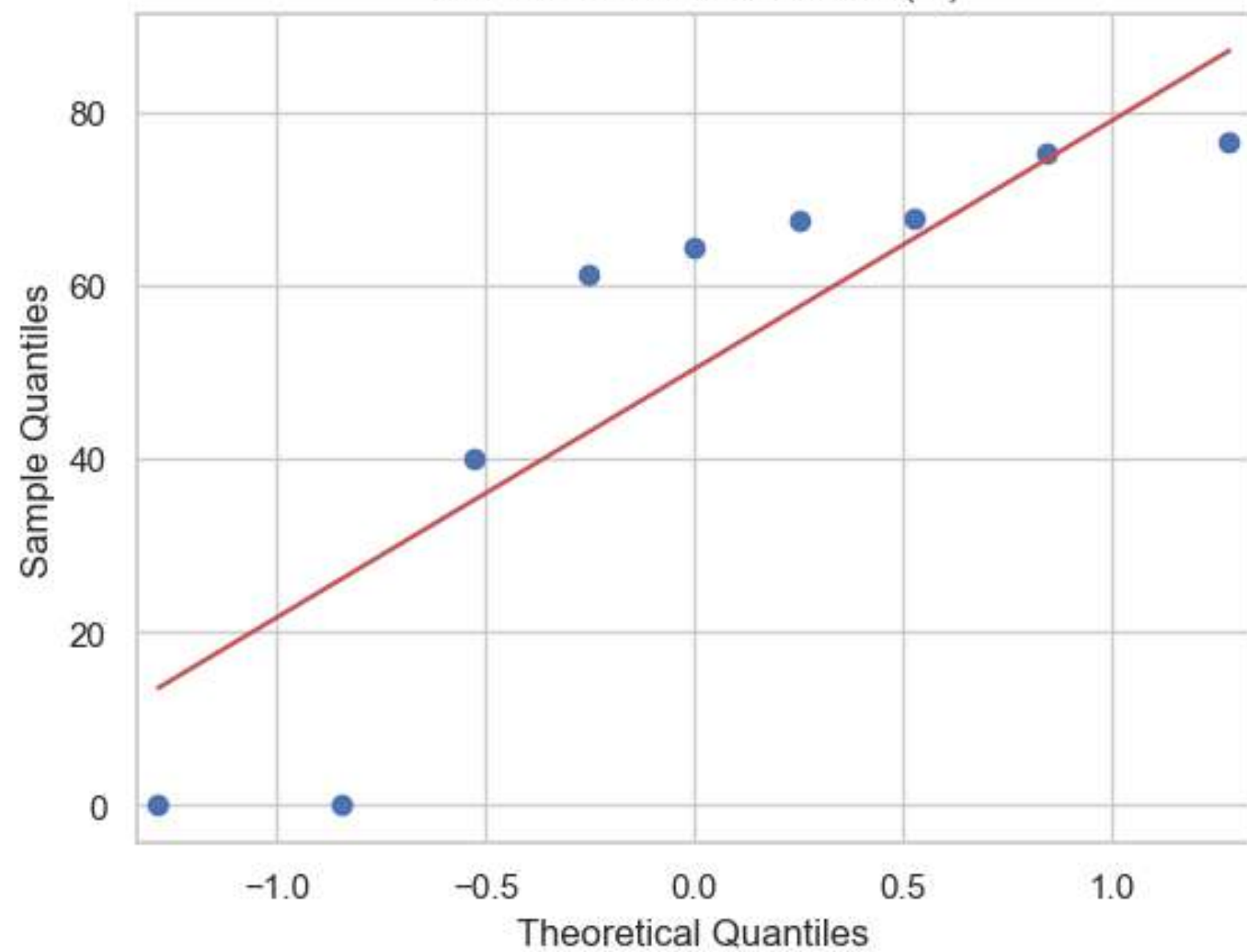


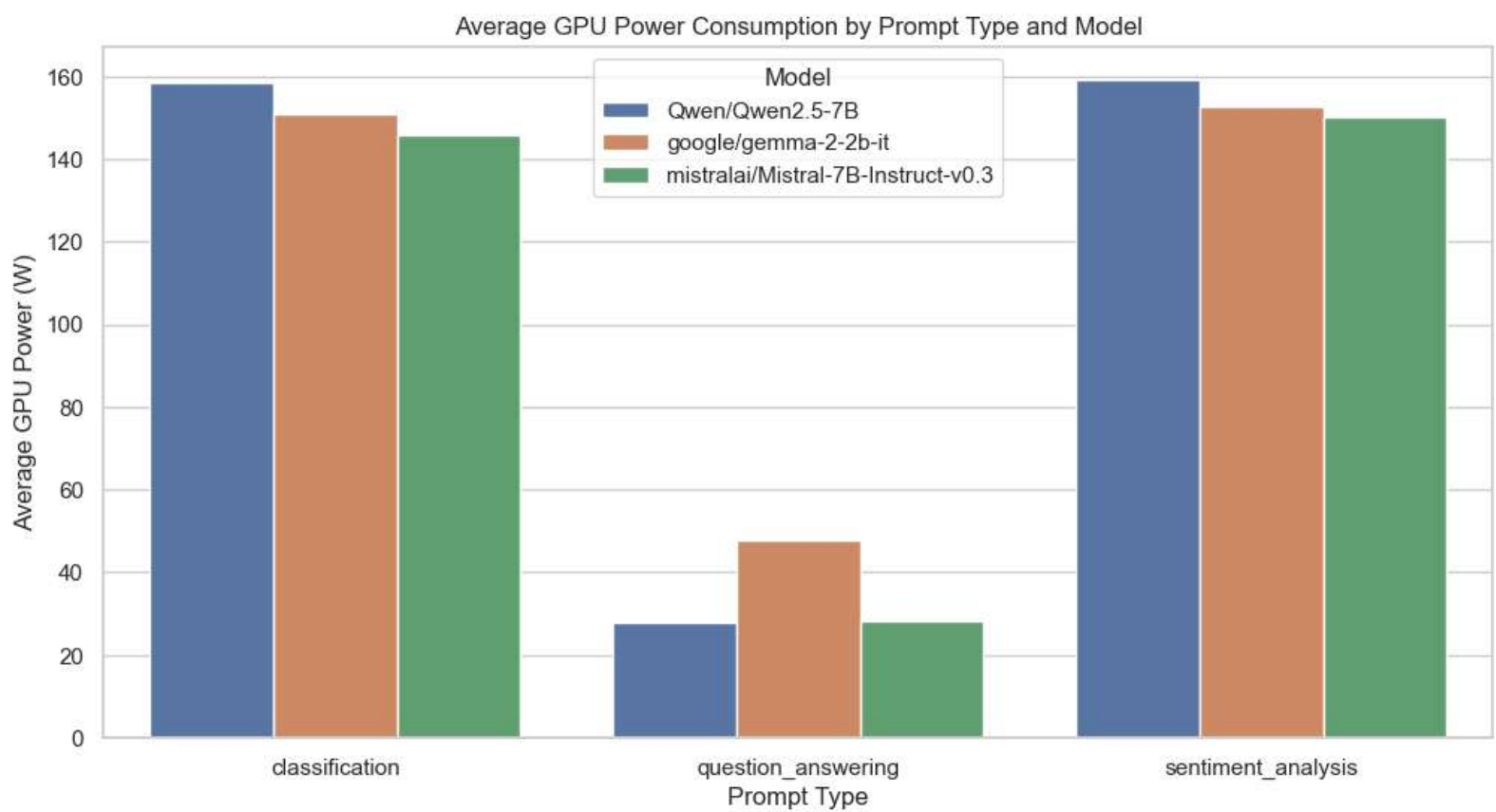


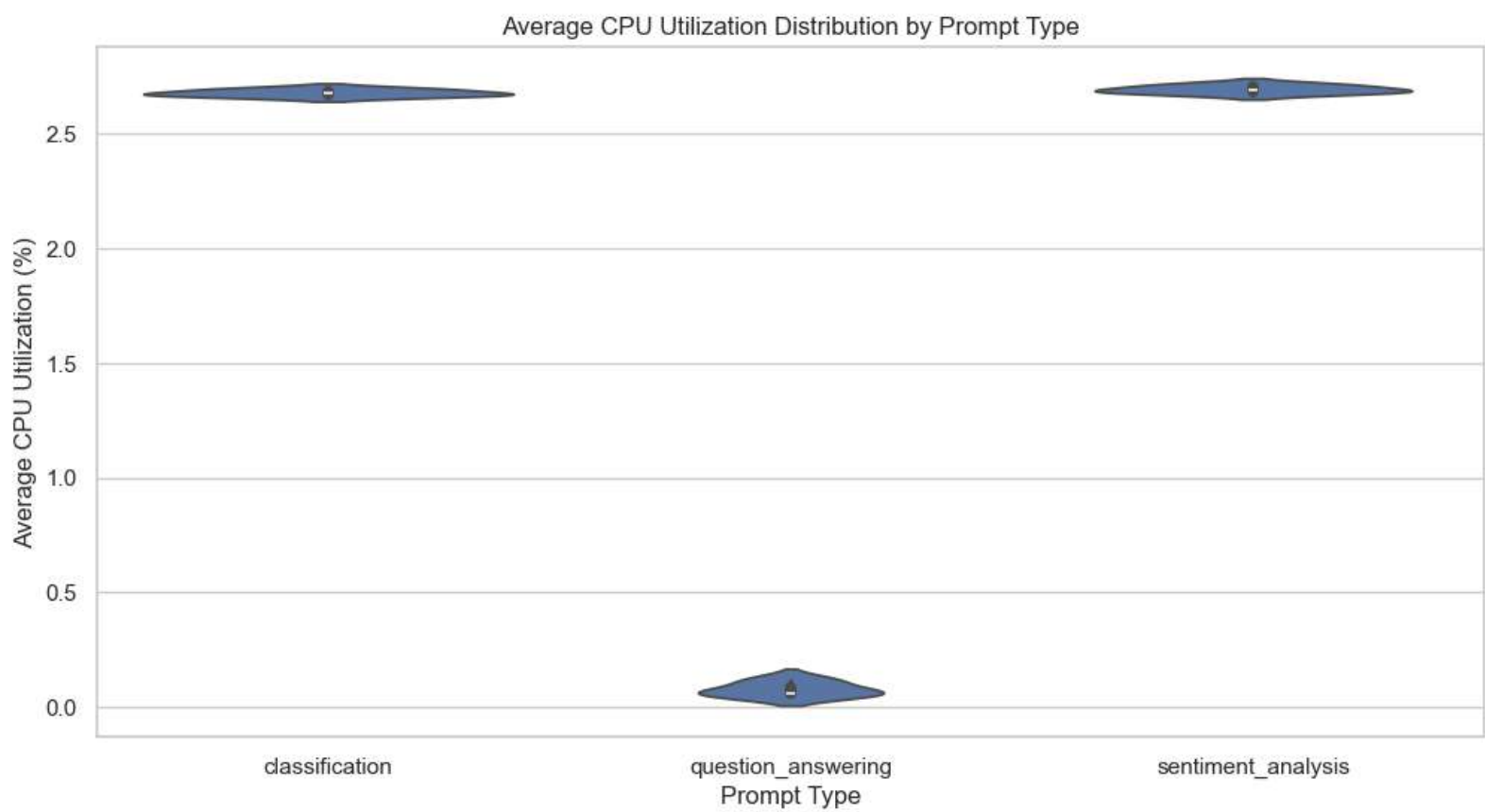
Q-Q Plot for Duration (s)



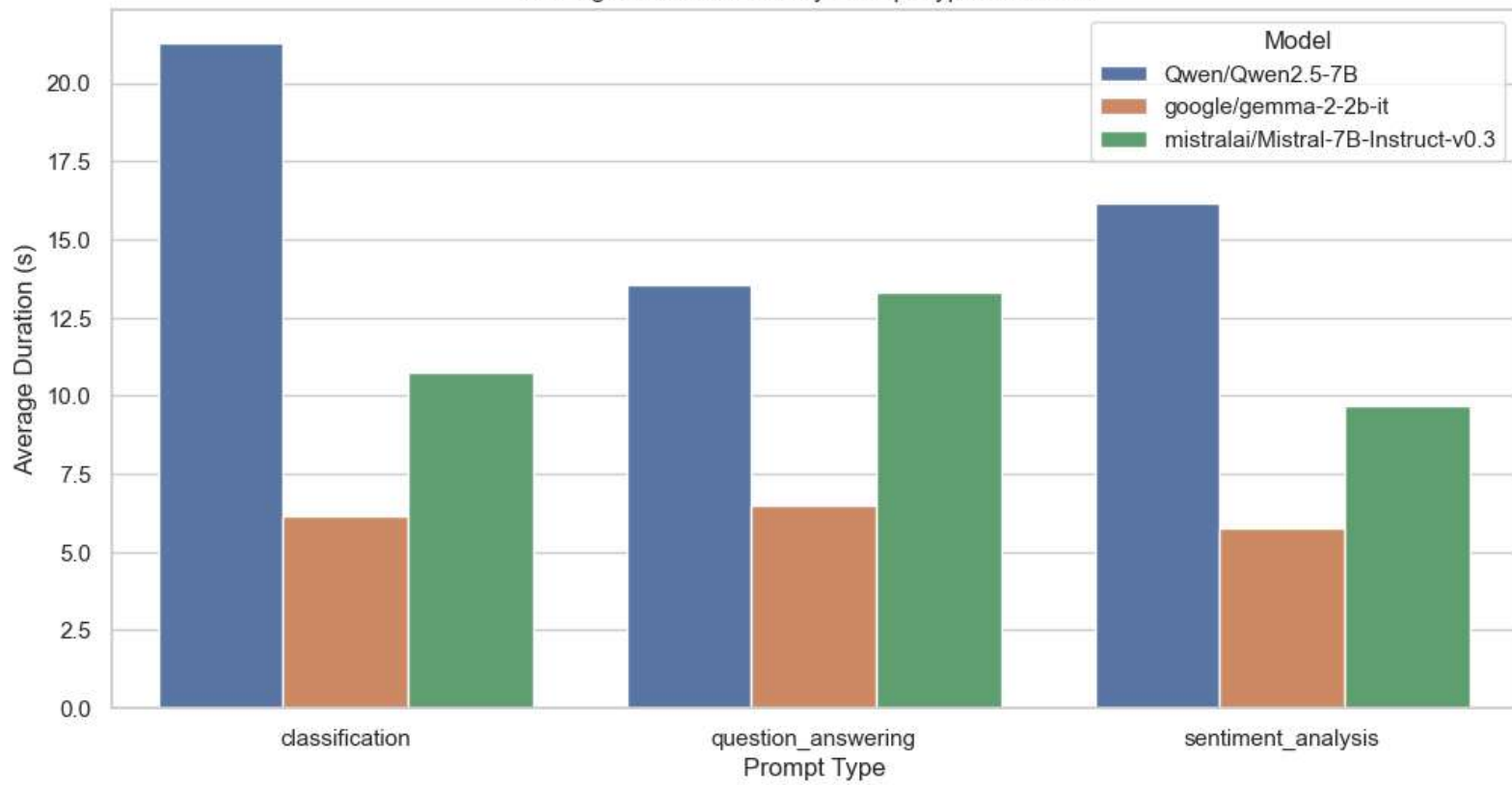
Q-Q Plot for GPU Utilization (%)



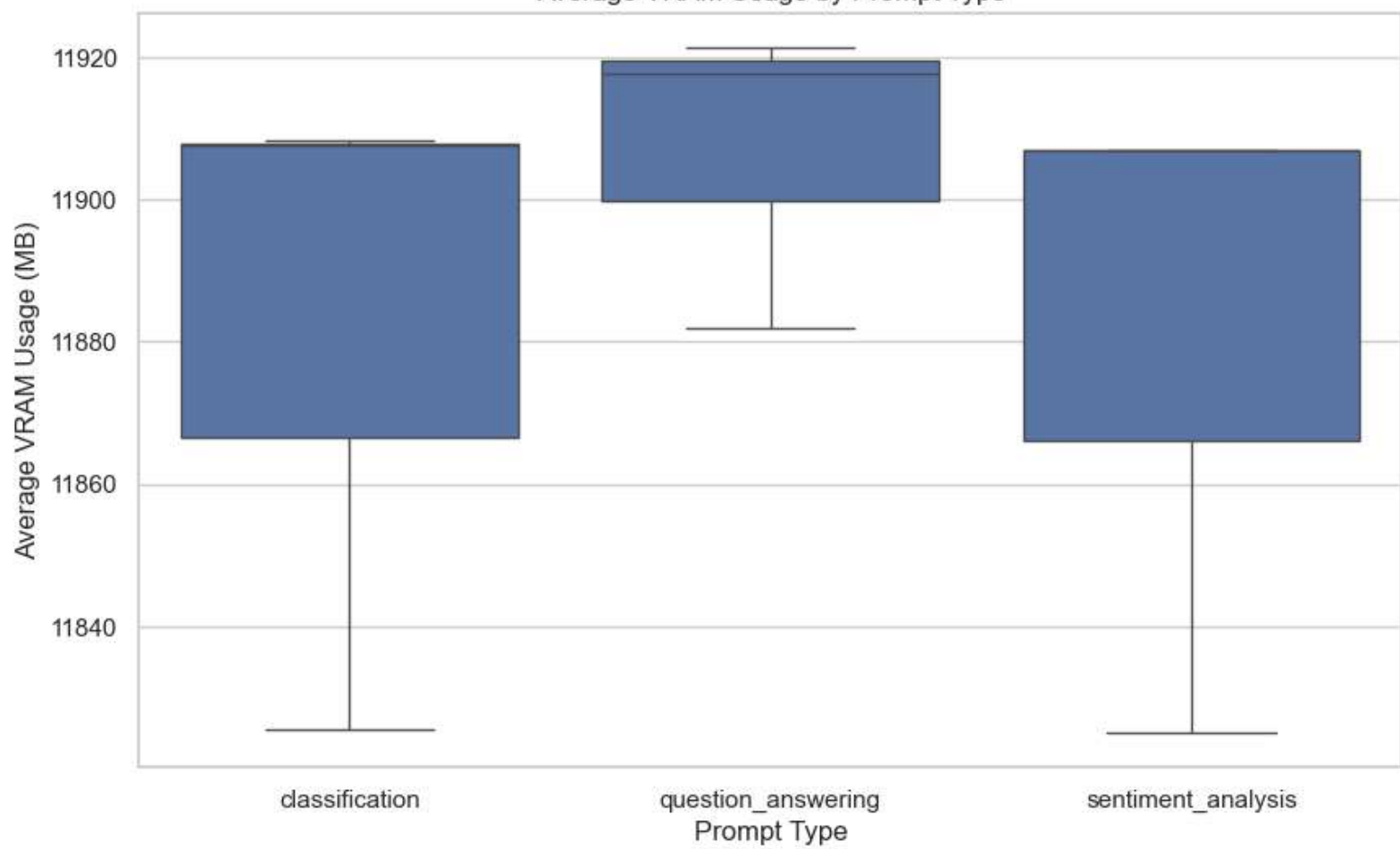


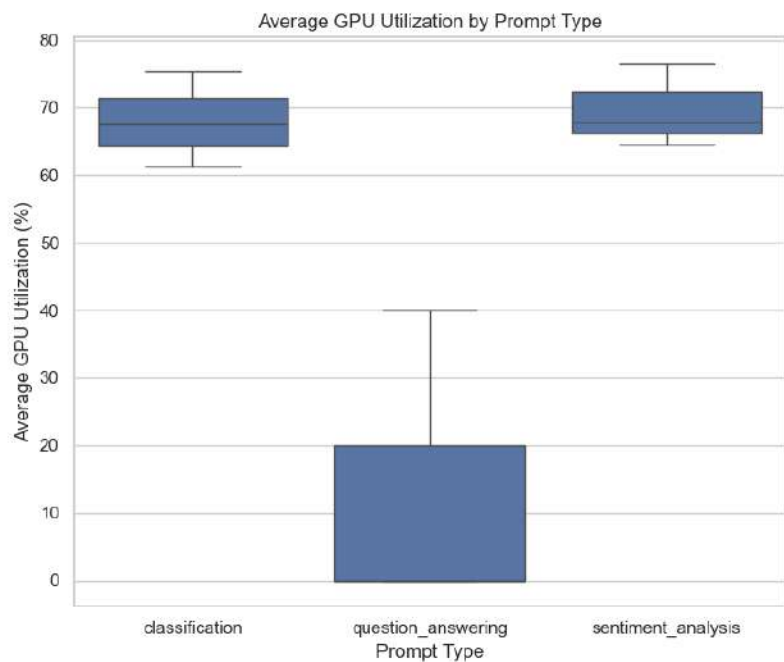
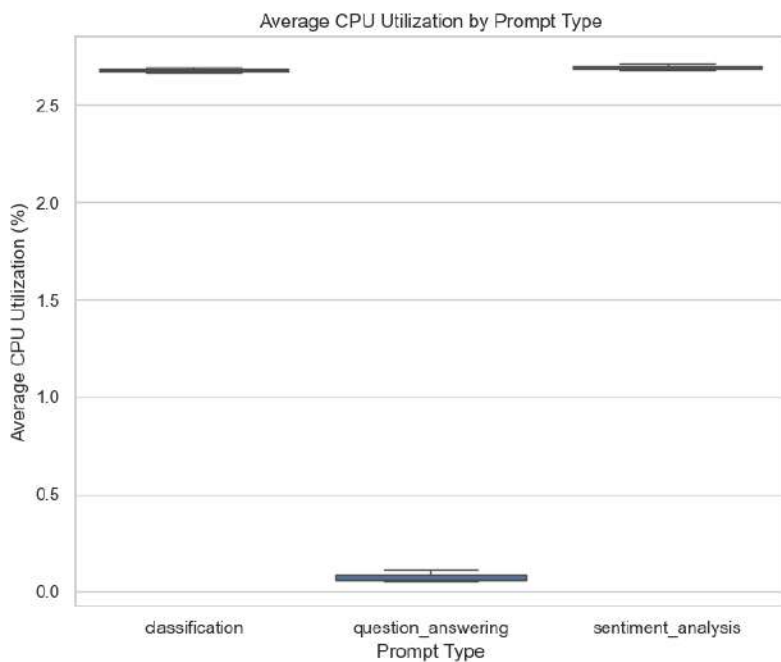


Average Inference Time by Prompt Type and Model

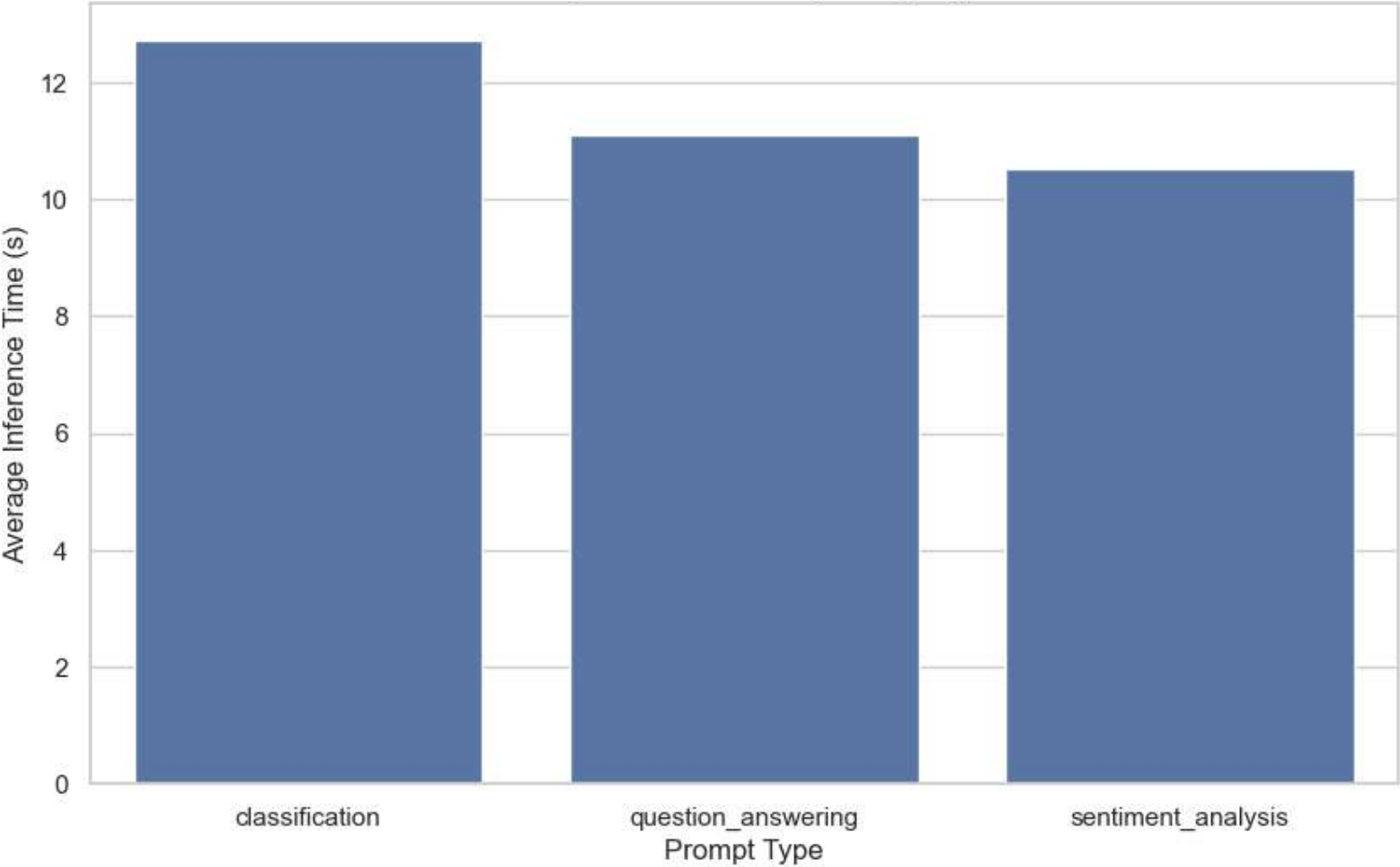


Average VRAM Usage by Prompt Type



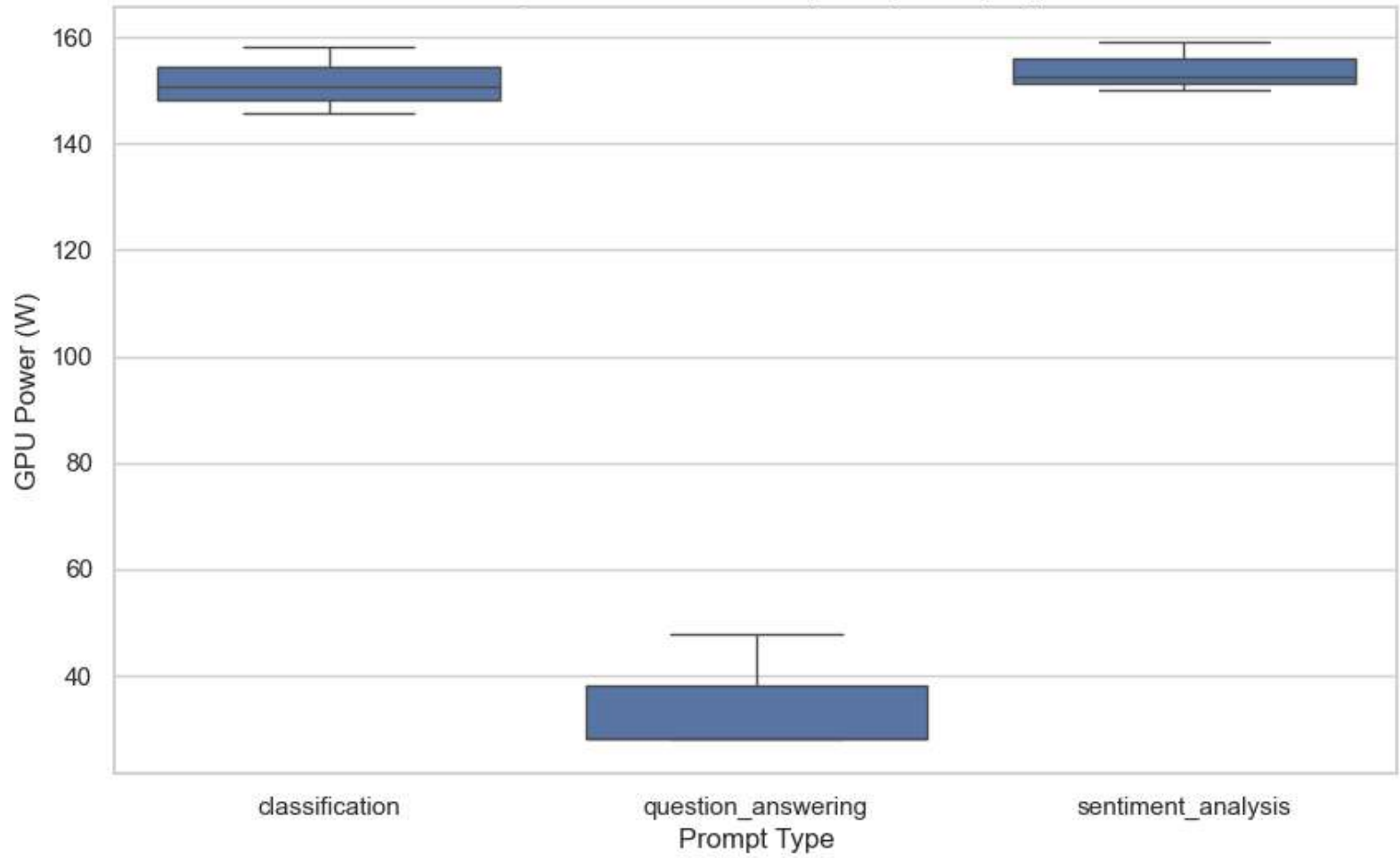


Average Inference Time by Prompt Type

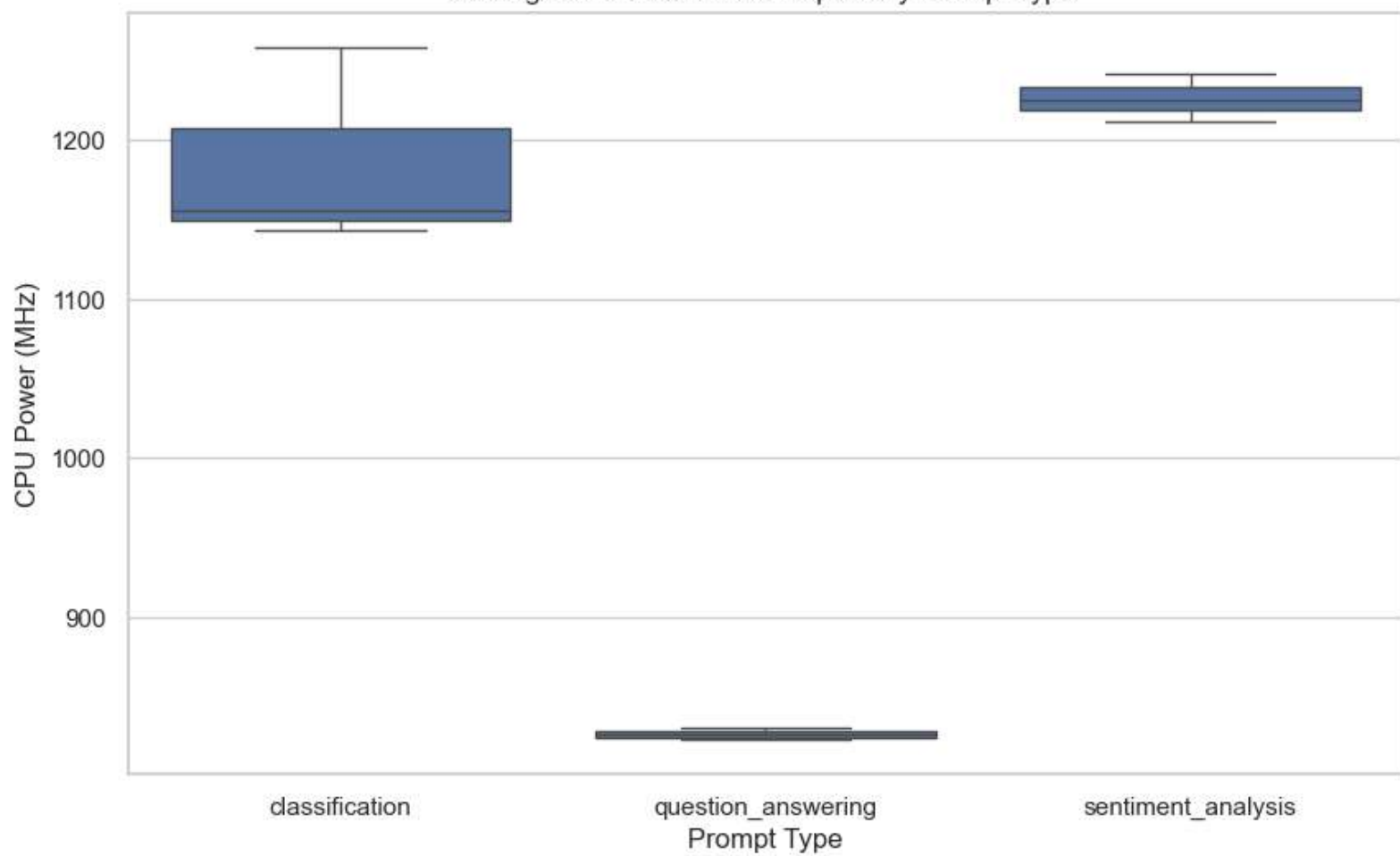




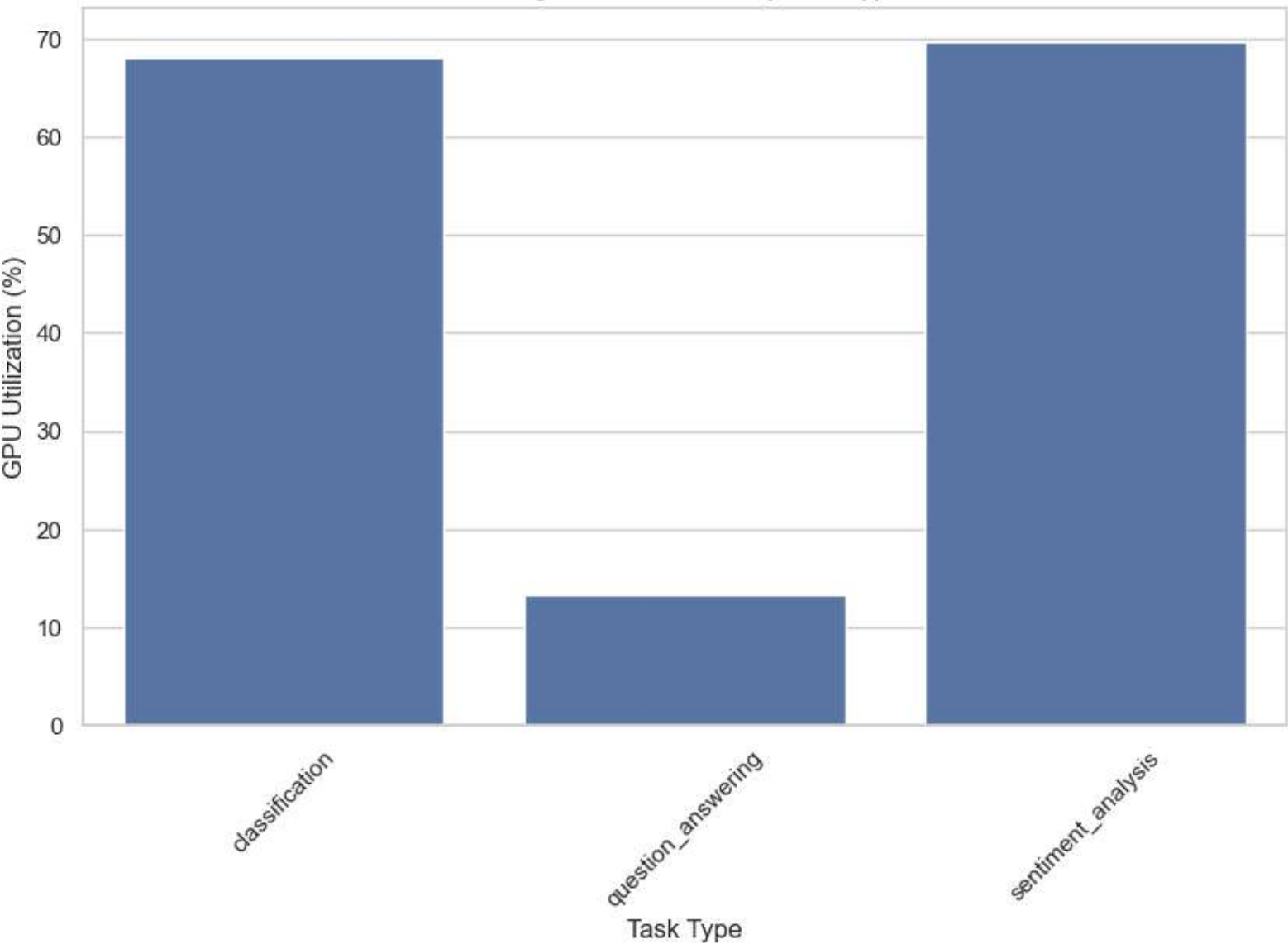
Average GPU Power Consumption by Prompt Type



Average CPU Power Consumption by Prompt Type



Average GPU Utilization by Task Type



Average GPU Utilization by Model

