**Model Performance vs. Compression Trade-off for GPT-J Edge-Cloud Collaboration**