# Deploying AI Models with Speed, Efficiency and Versatility
## *Inference on NVIDIA's AI Platform*

Whitepaper