# Challenges and Solutions for Efficient LLM Deployment

PD Dr. Haojin Yang,

Multimedia and Machine Learning Group,

Hasso Plattner Institute

**Design IT.**
**Create Knowledge.**

www.hpi.de

# The Evolution of Large Language Models

- From Transformer architecture (2017) to GPT-3 (2023): rapid technological evolution

- Scale explosion: Models growing from billions to trillions of parameters

- Breakthrough in natural language understanding and generation capabilities

- Fundamental shift from task-specific models to versatile foundation models

- Computational challenges: increasing resource requirements for deployment

- Transformative impact across industries and human-AI interaction paradigms

# Key Deployment Challenges

- High computational requirements for model inference

- Memory constraints with larger model architectures

- Latency issues affecting real-time applications

- Scaling difficulties across distributed systems

- Cost management for cloud-based deployments

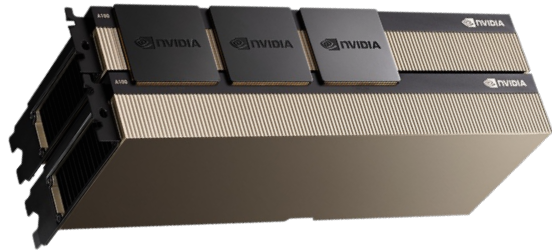- Integration challenges with existing infrastructure



Image source: https://hatchworks.com/blog/gen-ai/how-to-deploy-llm/
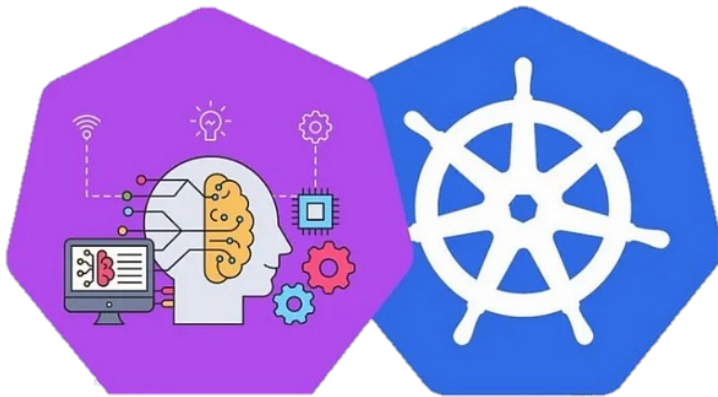
# Optimization Strategies

- Quantization: Reducing numerical precision without significant accuracy loss

- Knowledge Distillation: Creating smaller models that mimic larger ones

- Model Pruning: Removing redundant parameters and connections

- Efficient Architecture Design: Using techniques like MoE (Mixture of Experts)

- Caching and Inference Optimization: Storing common outputs to reduce computation

- Hardware-specific Acceleration: Leveraging specialized chips like GPUs and TPUs

23.04.25

Image source: Nvidia

# Architectural Solutions

- Distributed inference architecture for horizontal scaling

- Model quantization to reduce memory footprint

- Serverless deployment for on-demand scaling

- Containerization with orchestration (Kubernetes)

- Microservices architecture for modular LLM components

- API gateways and load balancing for traffic management

Image source: https://medium.com/@bijit211987/kubernetes-is-platform-of-choice-for-deploying-llms-cd12be2cfe59

Image source: Nvidia

# Future Directions & Best Practices

- Adopt model quantization techniques to reduce computational footprint

- Implement caching strategies for repeated queries and patterns

- Establish robust monitoring and observability practices

- Design for graceful degradation during peak loads

- Utilize hybrid deployment approaches (edge + cloud)

- Create comprehensive evaluation frameworks for continuous optimization



Image source: https://quantumzeitgeist.com/