



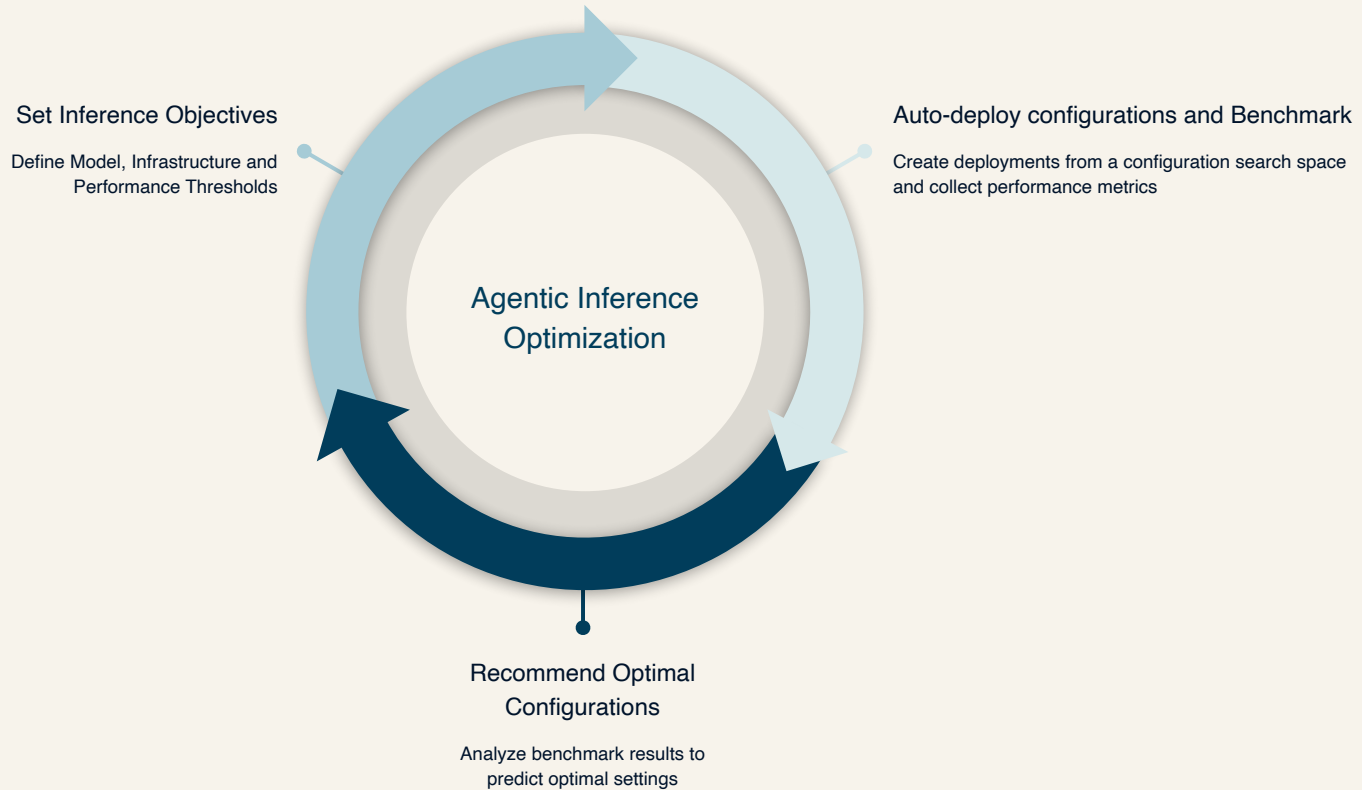
---

# Agentic Solutions For Inference Optimization

2026 MLSys Young Professionals Symposium

Kel Vanee, MVP, Machine Learning Engineering

# Agentic inference optimization



# Scale Engine Tuning Through Automation

## AI-Driven System Optimization

Automated Inference Performance Tuning

[Configuration](#) [Optimization](#) [Results](#)

### MODEL

Model Preset:

Name:  Served Name:

HF ID:  Tensor Parallel Size:

Model Path:

### PERFORMANCE TEST

Model Name:  Input Lengths:

Output Lengths:  Max Concurrency:

Request Rate (RPS):  Num Prompts:

Max Confgs:

### FIXED SERVER ARGS (APPLIED TO EVERY CONFIG)

[Add Arg](#)

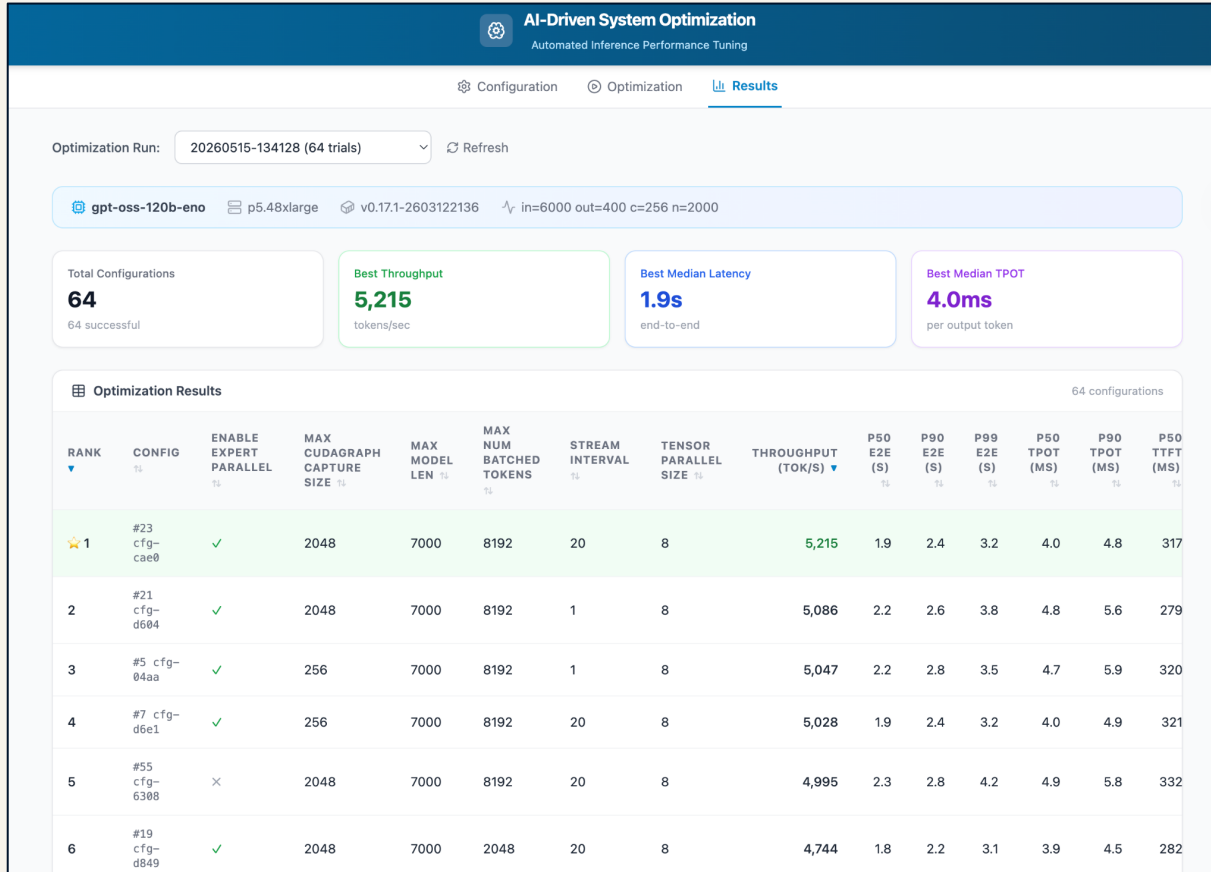
### SEARCH SPACE

Options (comma-separated):

Options (comma-separated):

Options (comma-separated):

# Scale Engine Tuning Through Automation



Note: the above represents sample data and is provided for illustrative purposes only

# Capital One AI Research Internships

We can't do it without great talent like you. Join us.

*WANT TO LEARN MORE?  
MEET US AT BOOTH #8*

## Applied Research Internship Program (ARIP)

Ph.D. students can expand their AI knowledge and move their research forward through Capital One's Applied Research Internship Program. During the 12-week summer internship, you'll engage in high-impact applied research, exploring novel and interesting AI challenges that will help create transformative customer experiences.



## Data Science Internship Program (DSIP)

Tap into the latest computing and machine learning tech in this program for Master's and PhD students. During the 10-week summer internship, you'll create dynamic models and leverage a diverse tech stack as you work across teams to unlock opportunities that improve the lives of our customers.



## AI Engineering Internship Program

Master's and PhD students can be part of a team delivering industry-leading capabilities and scalable, high-performance AI infrastructure. During this 12-week program, you'll help build AI and Agentic AI applications and platform capabilities that power customer facing experiences, cutting-edge GenAI workflows, shared infrastructure, and internal developer tools.

