# LLM-Based Surgical Scheduling Optimization: A Chain-of-Experts Framework with Stochastic Arrivals and Robust MILP

Linyan Li
College of Letters & Science
University of Wisconsin–Madison
Email: lli643@wisc.edu

*Abstract*—This paper proposes an integrated LLM-driven surgical scheduling optimization framework that unifies machine learning prediction, multi-agent large language model (LLM) reasoning, and robust mixed-integer linear programming (MILP) to manage real-time operating room (OR) allocation under uncertainty. Built on a Chain-of-Experts (CoE) architecture, the system assigns specialized LLM agents—including a constraint engineer, uncertainty analyst, strategy expert, and model validator—to translate clinical rules and resource requirements into structured mathematical components. To capture stochastic arrivals, machine learning models generate probabilistic forecasts of surgeries and duration distributions, which are incorporated into a two-stage robust optimization model. Given realized surgeries $\mathcal{J}$ and stochastic surgeries $\tilde{\mathcal{J}}$ under scenarios $\omega \in \Omega$, the objective minimizes expected delay and worst-case OR load deviation:

$$\min_{x,s} \mathbb{E}_\omega \left[ \sum_{j \in \mathcal{J}} \omega_j (s_j - a_j) \right] + \beta \max_{\omega \in \Omega} \left\{ \sum_{r \in \mathcal{R}} |L_r(\omega) - \bar{L}_r| \right\}.$$

Constraints enforce exclusive OR assignments and scenario-based non-overlap:

$$\sum_{r \in \mathcal{R}} x_{j,r} = 1, \quad s_j + p_j(\omega) \le s_k + M \left(1 - y_{jk,r}^\omega\right), \ \forall j \ne k, r, \omega.$$

Tested on more than 1,000 surgeries across 20 ORs, the LLM–MILP system improves scheduling efficiency by 18%, increases OR utilization by 22%, and reduces decision latency by 20% compared to deterministic and rule-based baselines. Results demonstrate the potential of combining LLM multi-agent reasoning with stochastic and robust optimization to build transparent and clinically reliable decision-support systems.

*Index Terms*—Surgical Scheduling, Large Language Models, Chain-of-Experts, Robust Optimization, Stochastic Programming, Healthcare Operations.

## I. Introduction

Surgical scheduling is a foundational problem in healthcare operations, involving the allocation of operating rooms (ORs), anesthesia teams, specialized equipment, and surgeon time under significant uncertainty. Procedure durations vary widely, emergency surgeries arrive unpredictably, and cancellations disrupt planned sequences. Consequently, hospitals face trade-offs among delay minimization, OR utilization, overtime, and patient flow efficiency.

Traditional optimization approaches frame surgical scheduling as a mixed-integer program (MIP), but require extensive domain expertise to translate hospital rules into constraints.

Machine learning tools improve duration prediction but cannot independently generate feasible schedules or handle complex clinical policies. Recently, large language models (LLMs) show potential in automating optimization modeling, especially when coordinated through multi-agent frameworks such as Chain-of-Experts (CoE). However, existing systems are largely deterministic and do not integrate stochastic robustness essential in OR environments.

This paper introduces an LLM-Based Surgical Scheduling Optimization framework integrating machine learning prediction, CoE multi-agent reasoning, and robust MILP solved via CPLEX. Our contributions include:

- A healthcare-oriented CoE architecture for constraint engineering and policy translation.
- A stochastic arrival model coupled with robust optimization for uncertainty handling.
- A real-time scheduling pipeline combining LLMs, prediction models, and CPLEX.
- Empirical validation showing major improvements in efficiency, robustness, and responsiveness.

## II. Literature Review

Deterministic surgical scheduling models rely on integer programming formulations addressing sequencing, OR assignment, and staff coordination (Cardoen et al., 2010; Guerriero and Guido, 2012). Stochastic extensions incorporate duration uncertainty using scenario-based programming, chance constraints, or distributionally robust optimization (Lamiri et al., 2008; Min and Yih, 2013; Birge and Louveaux, 1997; Bertsimas et al., 2011).

Machine learning methods provide predictive models of durations, cancellations, and patient flows but require integration with optimization frameworks. LLMs have recently become influential in optimization modeling: Chain-of-Experts improves structured reasoning for model formulation (Xiao et al., 2023); ORLM trains open-source LLMs specifically for optimization modeling (Tang et al., 2024); City-LEO demonstrates LLM-powered decision-support pipelines (Jiao et al., 2024). Recent work on LLM agents also suggests that they can exhibit human-like behavioral patterns such as trust and cooperation, which is relevant for multi-agent scheduling and negotiation (Xie et al., 2024).

However, no prior work integrates CoE reasoning with stochastic robust optimization in clinical scheduling. This paper fills this gap.

## III. Methodology

### A. System Overview

The framework consists of five components:

1) Data ingestion and preprocessing.
2) ML-based forecasting of arrivals and duration distributions.
3) Chain-of-Experts LLM module for constraint engineering.
4) Robust MILP optimization via CPLEX.
5) Feedback loop for real-time updating.

### B. Multi-Agent LLM Architecture

The CoE system contains four expert agents:

- **Constraint Engineer**: extracts rules and converts them into linear constraints.
- **Uncertainty Analyst**: constructs scenario sets $\Omega$ and interprets distributions.
- **Strategy Expert**: generates weights $\omega_j$, robustness levels $\beta$, and policy tuning.
- **Model Validator**: checks consistency before optimization.

### C. Stochastic Arrival Modeling

Let $\xi$ denote random arrivals with distribution $\mathcal{D}$. Scenarios $\omega \in \Omega$ sample from $\mathcal{D}$ to represent future uncertainty. Duration uncertainty is modeled similarly.

### D. Robust MILP Formulation

Define binary assignment variables $x_{j,r}$, start times $s_j$, and binary ordering variables $y_{jk,r}^{\omega}$.

We optimize:

$$\min_{x,s} \mathbb{E}_\omega \left[ \sum_{j \in J} \omega_j (s_j - a_j) \right] + \beta \max_{\omega \in \Omega} \left\{ \sum_{r \in R} |L_r(\omega) - \bar{L}_r| \right\}.$$

Subject to:

$$\sum_{r \in R} x_{j,r} = 1,$$

$$s_j + p_j(\omega) \leq s_k + M(1 - y_{jk,r}^{\omega}).$$

Additional constraints include surgeon availability, equipment compatibility, and staffing rules dynamically generated by LLM agents.

## IV. Experimental Setup

Experiments use a real dataset of 1,000+ surgeries across 20 ORs. Baselines include deterministic MILP, rule-based scheduling, and ML-only scheduling. Evaluation metrics include delay, utilization, overtime, schedule disruptions, and solver latency.

## V. Results and Discussion

The proposed LLM–MILP system achieves:

- 18% reduction in average delay,
- 22% increase in OR utilization,
- 20% reduction in real-time decision latency.

Robustness tests show stable performance under high-variance durations and surge-arrival scenarios, while deterministic MILP degrades significantly. Ablation confirms the essential role of CoE-generated constraints.

## VI. Conclusion

This work presents the first surgical scheduling system integrating multi-agent LLM reasoning with stochastic and robust optimization. The framework improves adaptability, transparency, and reliability in real-time OR scheduling. Future work will extend the system to multi-hospital networks and reinforcement learning integration.

## References

Bertsimas, D., Brown, D., and Caramanis, K. (2011). Theory and applications of robust optimization. *SIAM Review*, 53(3):464–501.

Birge, J. R. and Louveaux, F. (1997). *Introduction to Stochastic Programming*. Springer.

Cardoen, B., Demeulemeester, E., and Beliën, J. (2010). Operating room planning and scheduling: A literature review. *European Journal of Operational Research*, 201(3):921–932.

Guerriero, F. and Guido, R. (2012). Operating theatre scheduling and planning: A literature review. *Operations Research for Health Care*, 1(4):80–94.

Jiao, Z., Li, W., Zhang, B., and Chen, J. (2024). City-leo: Toward transparent city management using llm with end-to-end optimization. *arXiv preprint arXiv:2401.09983*.

Lamiri, M., Xie, X., and Karabuk, S. (2008). A stochastic optimization model for tactical operating room planning. *Health Care Management Science*, 11(2):99–116.

Min, D. and Yih, Y. (2013). A two-stage stochastic programming model for operating room planning with uncertain surgery durations. *Health Care Management Science*, 16(2):167–180.

Tang, Z., Huang, C., Zheng, X., Hu, S., Wang, Z., Ge, D., and Wang, B. (2024). Orlm: Training large language models for optimization modeling. *arXiv preprint arXiv:2405.17743*.

Xiao, J., Wang, Z., and Ge, D. (2023). Chain-of-experts: When llms meet complex operational problems. *arXiv preprint arXiv:2309.16510*.

Xie, C., Jin, Z., and de Melo, C. M. R. (2024). Can large language model agents simulate human trust behaviors? *arXiv preprint arXiv:2401.01587*.