

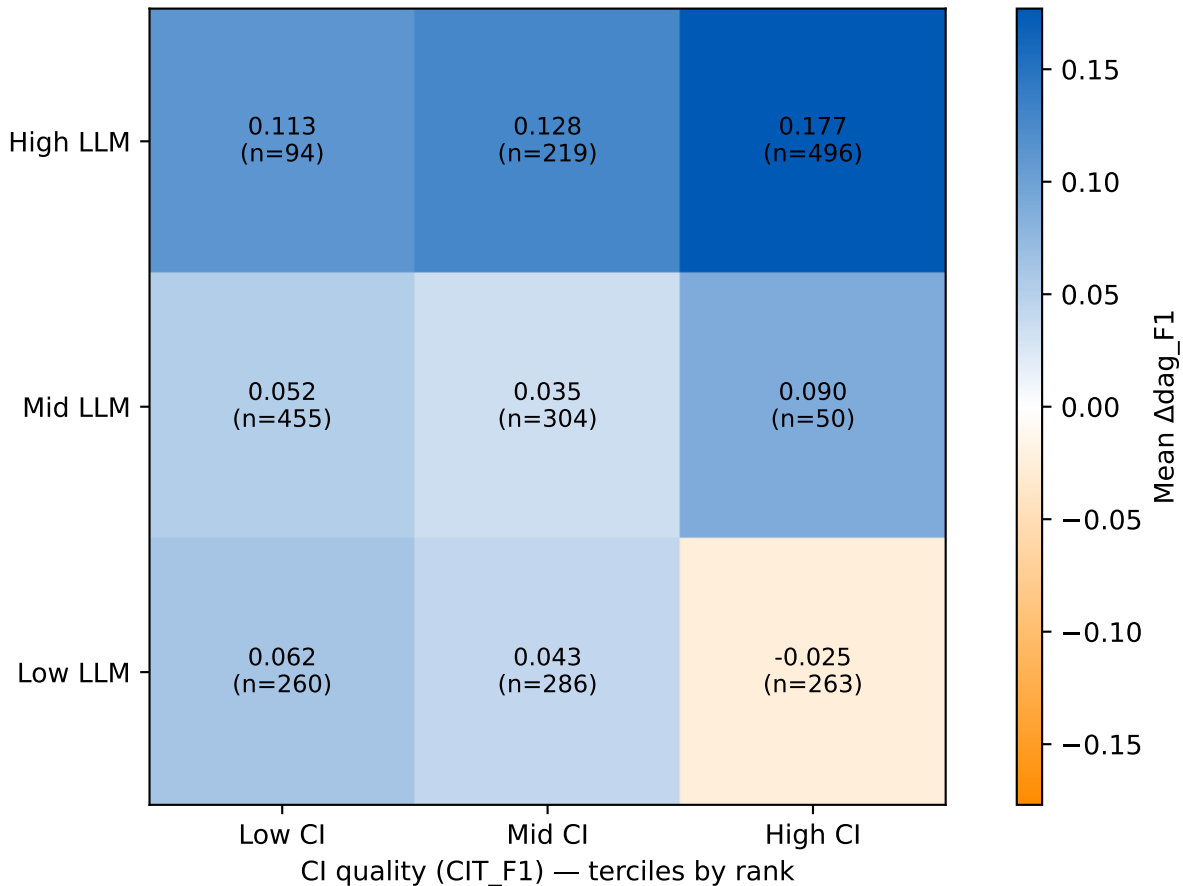
# Aggregated 3×3 Heatmaps — Rank-based Low/Mid/High

Each axis binned into terciles by rank (ensures ~equal counts per bin)

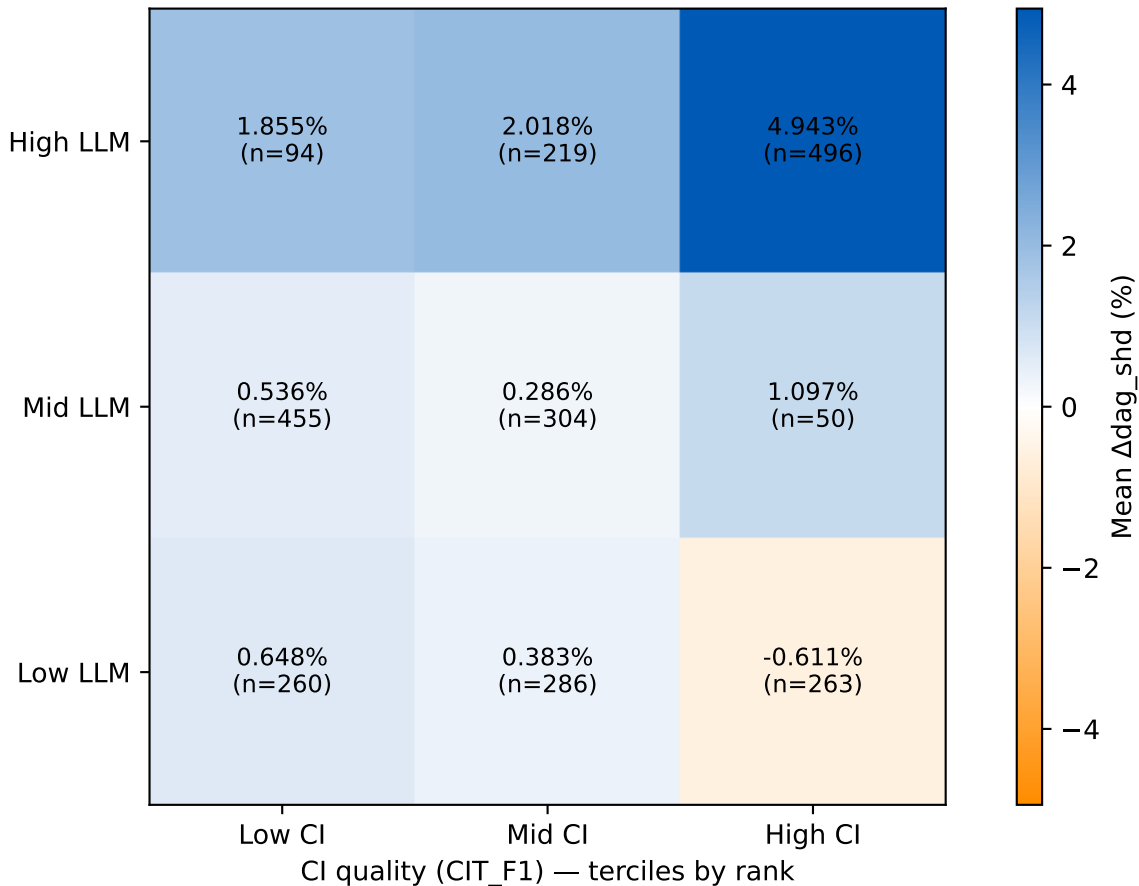
2025-10-02 12:50

# DAG Metrics

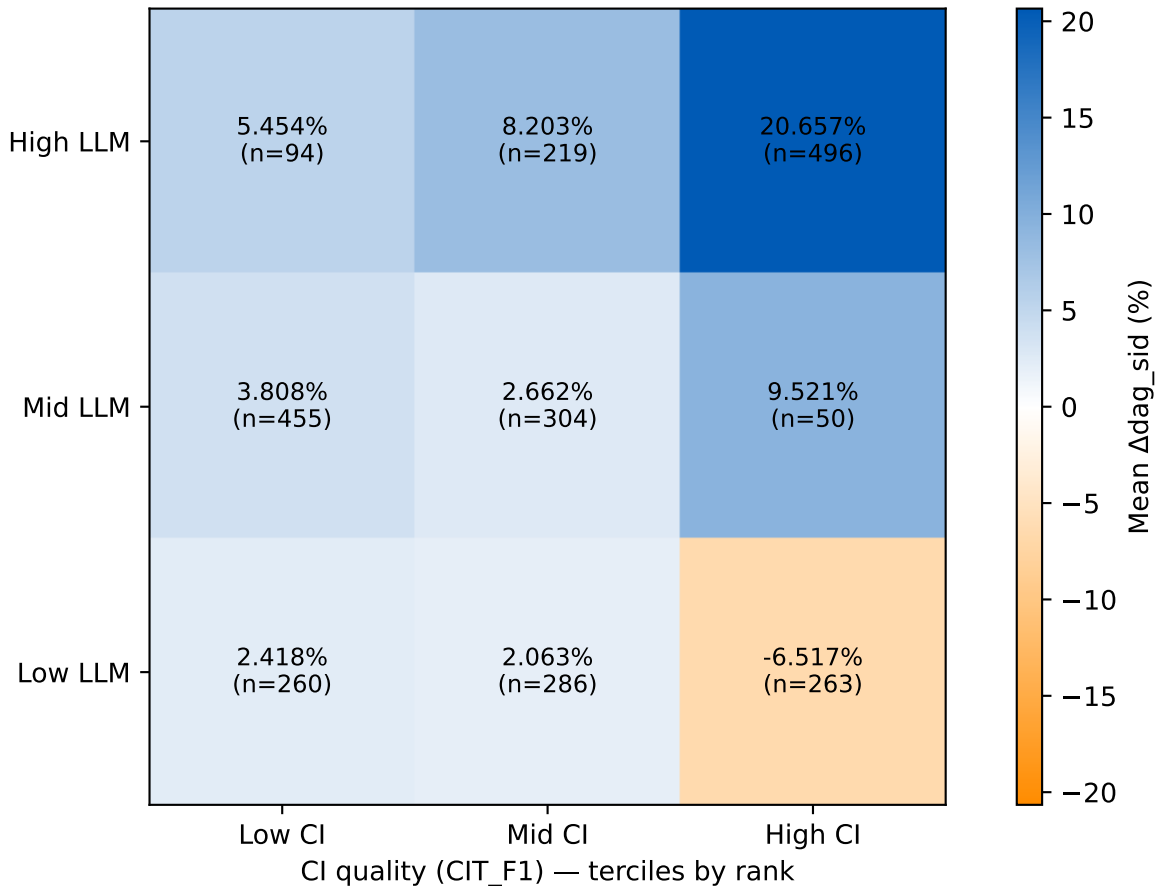
LLM precision — terciles by rank



LLM precision — terciles by rank



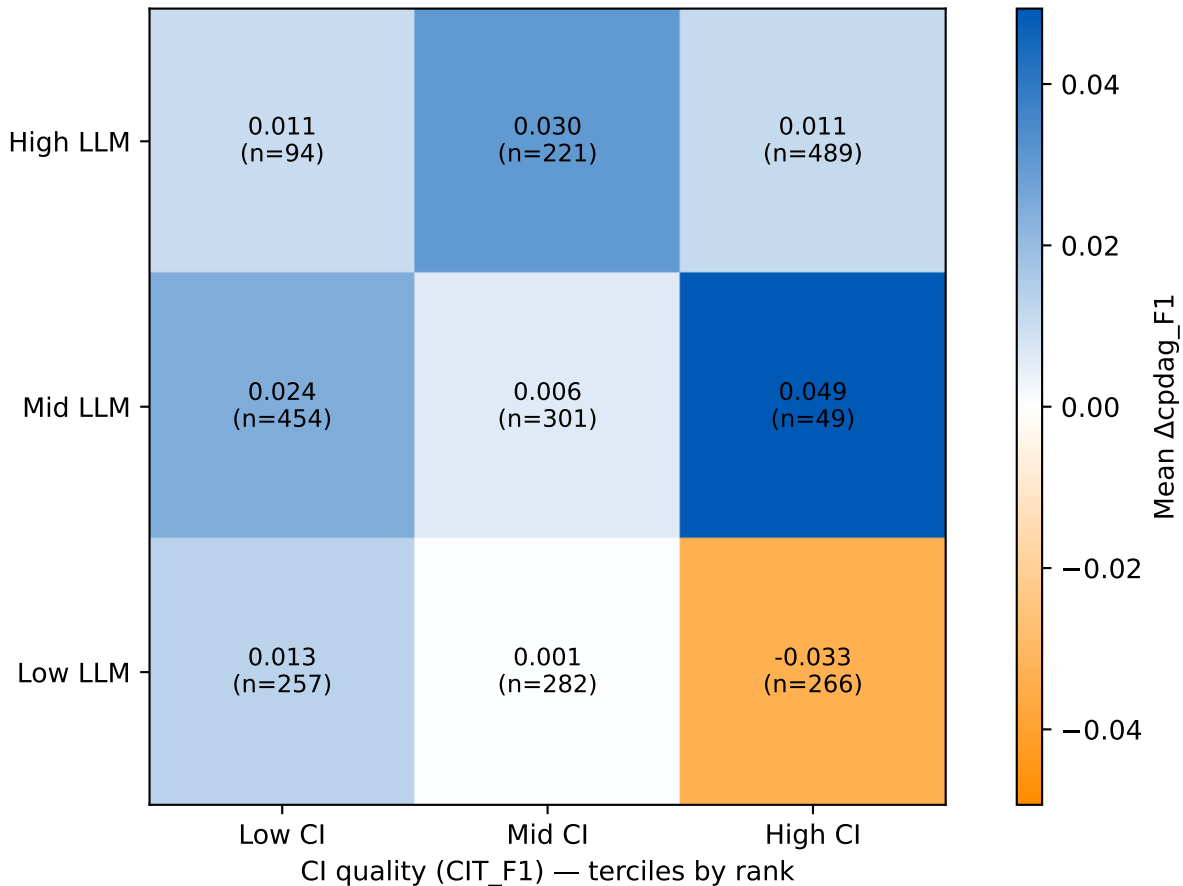
LLM precision — terciles by rank



# CPDAG Metrics

2025-10-02 12:50

LLM precision — terciles by rank



LLM precision — terciles by rank

High LLM

Mid LLM

Low LLM

0.214%  
(n=94)

0.519%  
(n=219)

0.783%  
(n=496)

0.440%  
(n=455)

0.046%  
(n=304)

2.590%  
(n=50)

0.316%  
(n=260)

-0.007%  
(n=286)

-4.036%  
(n=263)

Low CI

Mid CI

High CI

CI quality (CIT\_F1) — terciles by rank

Mean  $\Delta$ cpdag\_shd (%)

4

3

2

1

0

-1

-2

-3

-4



LLM precision — terciles by rank

High LLM

Mid LLM

Low LLM

0.792%  
(n=94)

4.399%  
(n=219)

4.384%  
(n=496)

3.527%  
(n=455)

0.781%  
(n=304)

13.191%  
(n=50)

1.440%  
(n=260)

0.792%  
(n=286)

-13.645%  
(n=263)

Low CI

Mid CI

High CI

CI quality (CIT\_F1) — terciles by rank

Mean  $\Delta$ cpdag\_sid\_avg (%)

10

5

0

-5

-10