

Compare Experiments

Compare evaluation results across multiple experiments side-by-side.

Navigation

[Dashboard](#)[Tables](#)[Ontology](#)[Labels](#)[Batches](#)

Annotations

[New Run](#)[Run History](#)

Evaluation

[New Evaluation](#)[Compare](#)

Comparison Setup

Select batch files and ground truth labels to compare.

Ground Truth Labels

[Ground Truth Small Table Based On Name](#)

Batch Files to Compare

[single | cot | BEO 431c litellm azure-gpt-4.1](#)[edm | cot | BEO 431c ollama gpt-oss:20b](#)[+ Add Another](#)[Compare Experiments](#)

Comparison Results

2 of 2 experiments compared successfully

Best Scores: Path F1: **41.45%**
Node F1: **34.84%**

Experiment	Provider	Model	Mode	Prompt	Path F1	Node F1	Columns	Status
batch_20251205_224711	litellm	azure-gpt-4.1	single	cot	<div style="width: 34.84%;">34.84%</div>	<div style="width: 41.45%;">41.45%</div>	431/431	Success
batch_20251205_224635	ollama	gpt-oss:20b	edm	cot	<div style="width: 23.51%;">23.51%</div>	<div style="width: 29.63%;">29.63%</div>	431/431	Success

batch_20251205_224711

litellm azure-gpt-4.1

single | cot

Path-Level

P: 37.43%

R: 34.57%

F1: **34.84%**

Node-Level

P: 46.33%

R: 40.26%

F1: **41.45%**

batch_20251205_224635

ollama gpt-oss:20b

edm | cot

Path-Level

P: 25.14%

R: 24.25%

F1: **23.51%**

Node-Level

P: 32.44%

R: 29.81%

F1: **29.63%**

[Settings](#)