Evaluation Report Visualization

generatedAtUtc: 2026-02-25T00:18:35.0450425+00:00
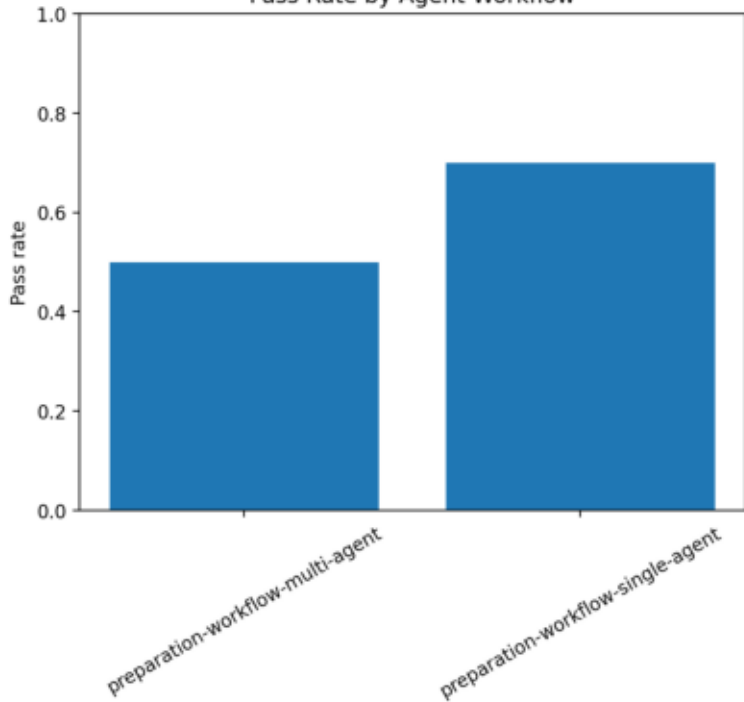totalCases: 20
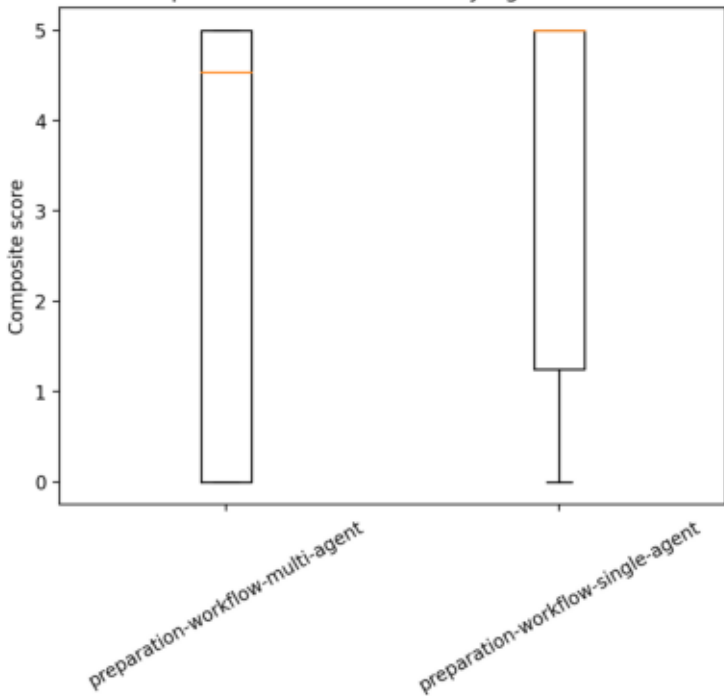totalPassed: 12
overallPassRate: 0.600

Agent summary:
- preparation-workflow-multi-agent: passRate=0.500 passed=5/10, compositeAvg(report)=2.908, metricCoverage=0.600
- preparation-workflow-single-agent: passRate=0.700 passed=7/10, compositeAvg(report)=3.000, metricCoverage=0.700
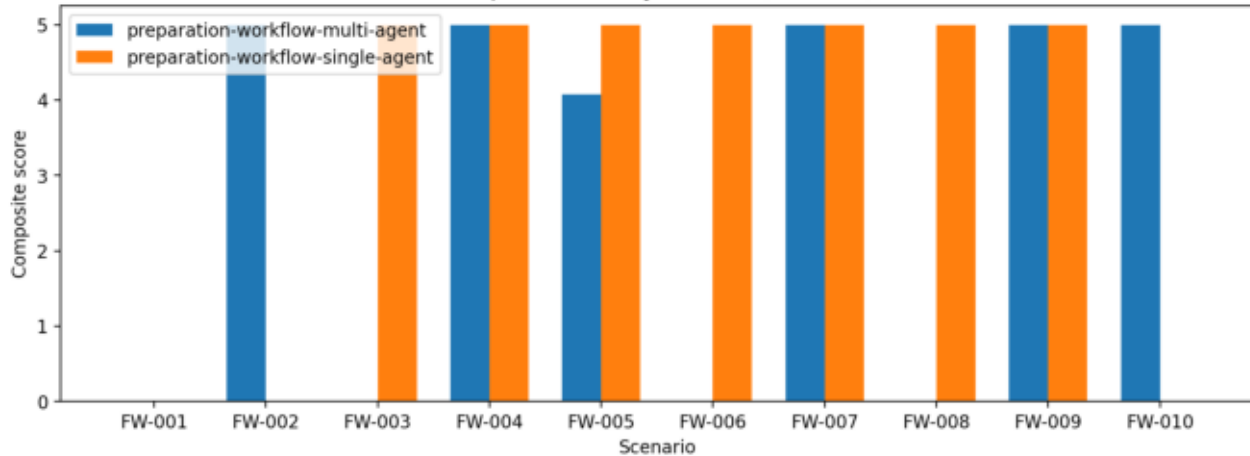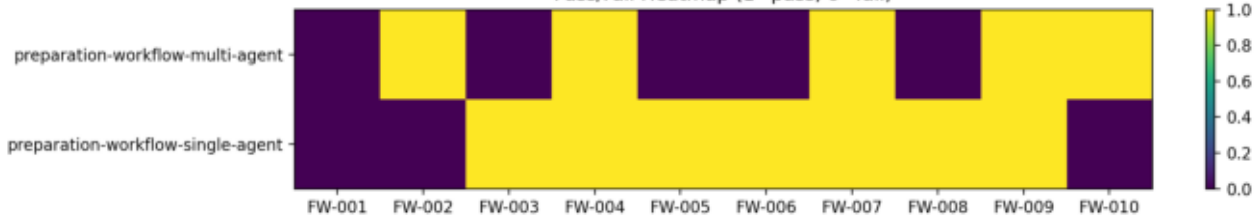
Pass Rate by Agent Workflow

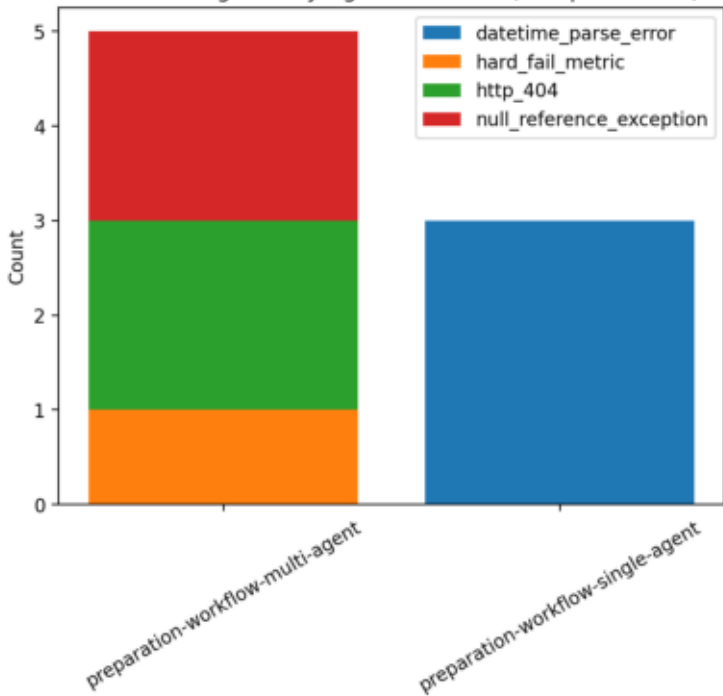Composite Score Distribution by Agent Workflow
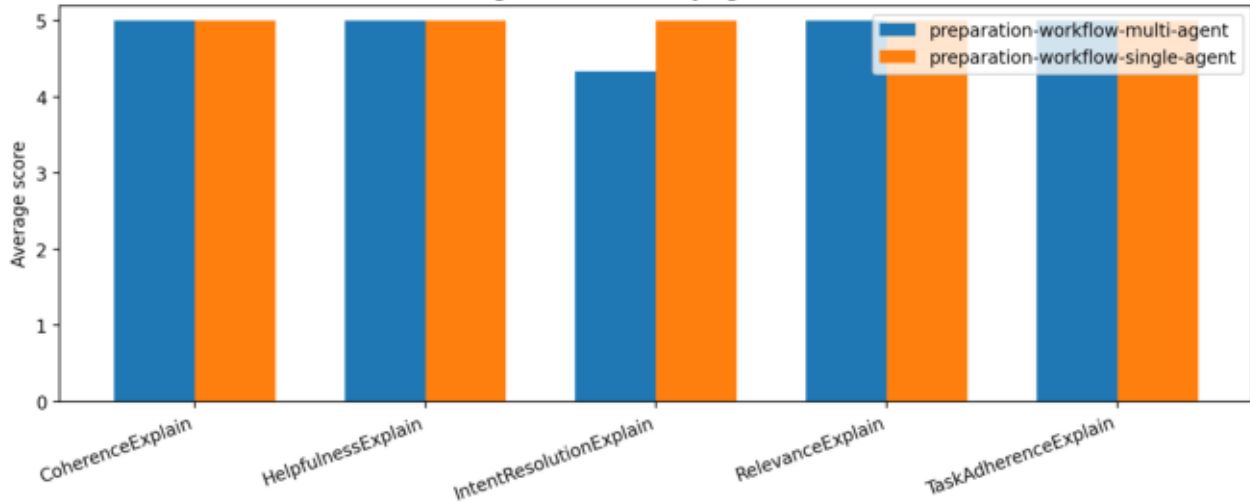
Composite Score by Scenario (FW-###)

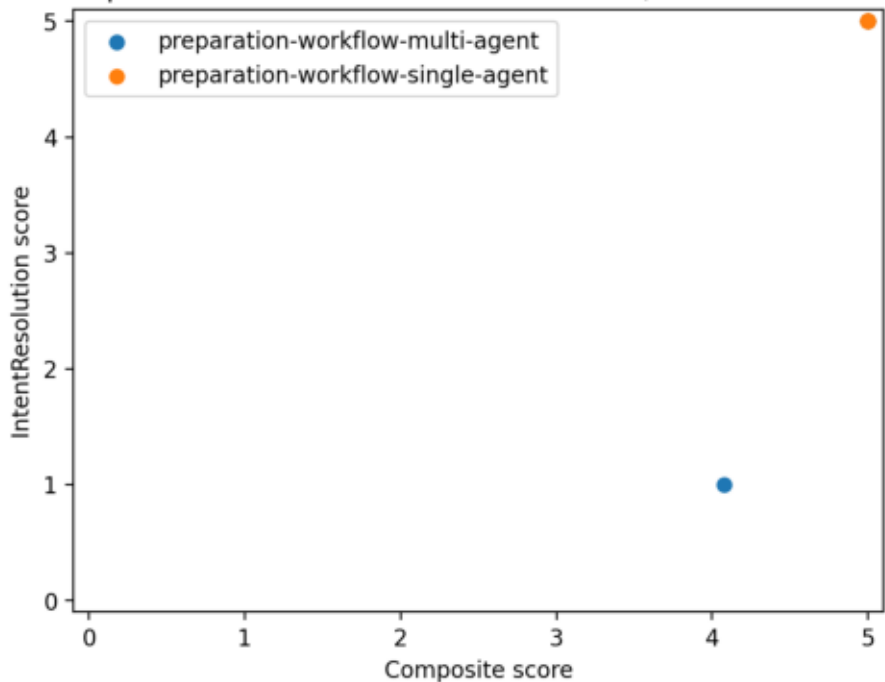Pass/Fail Heatmap (1=pass, 0=fail)

Failure Categories by Agent Workflow (non-pass cases)
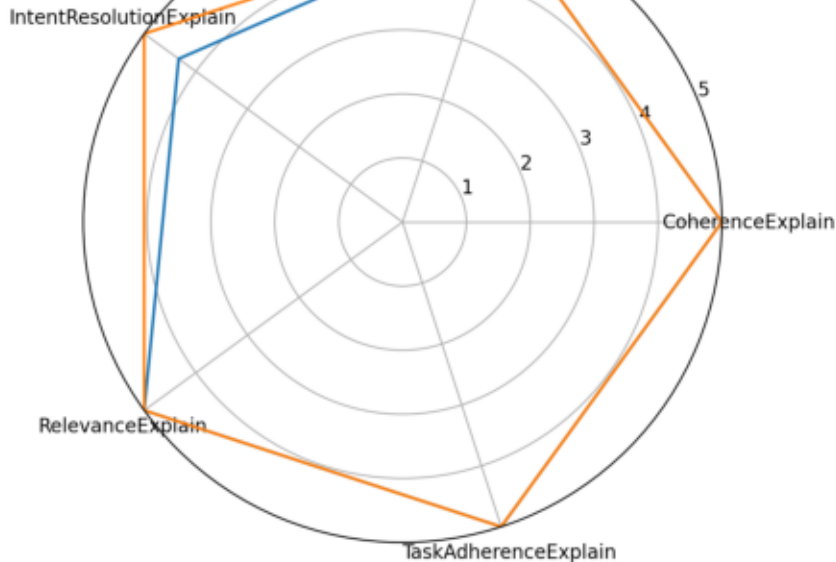
Average Metric Scores by Agent Workflow

Composite Score vs IntentResolution Score (cases with metrics)

Radar: Average Metric Scores

Legend:
- preparation-workflow-multi-agent
- preparation-workflow-single-agent

Axes: HelpfulnessExplain, CoherenceExplain, TaskAdherenceExplain, RelevanceExplain, IntentResolutionExplain

Metric-Evaluation Coverage by Workflow (cases with metrics / total)

# Scenario Comparison (Pivoted)

| scenario | compositeScore.preparation-workflow | compositeScore.preparation | failureCategory.preparation | failureCategory.preparation | passed.preparation-workflow | passed.preparation-workflow-single |
|---|---|---|---|---|---|---|
| FW-001 | 0.0 | 0.0 | null_reference_exception | datetime_parse_error | False | False |
| FW-002 | 5.0 | 0.0 | none | datetime_parse_error | True | False |
| FW-003 | 0.0 | 5.0 | http_404 | none | False | True |
| FW-004 | 5.0 | 5.0 | none | none | True | True |
| FW-005 | 4.076923076923077 | 5.0 | hard_fail_metric | none | False | True |
| FW-006 | 0.0 | 5.0 | null_reference_exception | none | False | True |
| FW-007 | 5.0 | 5.0 | none | none | True | True |
| FW-008 | 0.0 | 5.0 | http_404 | none | False | True |
| FW-009 | 5.0 | 5.0 | none | none | True | True |
| FW-010 | 5.0 | 0.0 | none | datetime_parse_error | True | False |