**grok - Phase 1: Comprehensive Metrics**
**(2 experiments average)**

| | Accuracy | Macro-F1 | Weighted-F1 |
|-------|----------|----------|-------------|
| Type | 0.700 | 0.727 | 0.704 |
| Party | 0.900 | 0.893 | 0.900 |
| Action | 0.683 | 0.454 | 0.696 |
| Asset | 0.600 | 0.146 | 0.606 |

Metric Score