

SPHERE: An Evaluation Card for Human-Al Systems sphere-eval.github.io/



Qianou Ma^{1*} (qianouma@cmu.edu), Dora Zhao^{2*} (dorothyz@stanford.edu), Xinran Zhao¹, Chenglei Si², Chenyang Yang¹, Ryan Louie², Ehud Reiter³, Diyi Yang²⁺, Tongshuang Wu¹⁺ (*Equal contribution, +Equal contribution)



¹Carnegie Mellon University,

²Stanford University,





Introducing a comprehensive template for designing and documenting human-Al system evaluations

Motivation

- Mounting concern over how to evaluate generative systems as existing methods fall short
- Human-Al systems introduce additional complexities as we must consider not only model performance but also the impact on users
- We introduce **SPHERE cards**, which provide a holistic overview on how to design and document the evaluations of human-Al systems

SPHERE (Subject-Process-Handler-Elapsed-Robustness Evaluation) Cards

What is being evaluated?

What components and system design goals are evaluated?

HOW is the evaluation conducted?

What is the <u>scope</u> of evaluation? What <u>methods</u> are used?

Who participates in the evaluation?

What are the <u>automated</u> or <u>human</u> evaluators used?

When is the evaluation conducted?

What is the time scale over which the evaluation is conducted?

HOW is the evaluation validated?

How do we ensure the <u>validation</u> of the evaluation?

Recommendations

Evaluations should reflect real—world use

Test systems in the real–world with evaluators from relevant stakeholder groups.

2) Cross-verify results across methods

Triangulate results across from different data sources or methods and ground methods in existing practices

Rigorously evaluate the evaluation

Expand practices for measuring reliability / validity and document practices for replication

Example SPHERE Card: AngleKindling (Petridis et al. 2023) What is being evaluated? The system's effectiveness (# of pursuable angles), efficiency (mental demand), and satisfaction (perceived helpfulness) How is the evaluation conducted? Extrinsic within-subject study with a quantitative post-task survey and *qualitative* interviews Who is participating in the evaluation? 12 professional journalists (domain experts, intended users) When is the evaluation conducted Short-term sessions of up to 60 minutes How is the evaluation validated? Counterbalancing of tool order to reduce learning effects for *validity*

