



From Data to Insights: Validating Artificial Intelligence (AI)-Generated Writing in Evidence Synthesis

Allie Cichewicz, MSc

¹Independent Consultant, Boston, MA, USA

Introduction

- Generative AI is increasingly integrated into systematic review software, progressing from structured tasks like screening and data extraction toward generating written summaries.
- Automated written synthesis remains a key unmet capability, and this shift to narrative output introduces new risks around accuracy, completeness, and misrepresentation of source findings that require formal validation methods.
- The RAISE (Responsible AI-Supported Evidence synthesis) framework provides guidance on transparent reporting of AI use in evidence synthesis, but practical criteria for evaluating the quality of AI-generated text against source literature are still lacking.

Objective

To evaluate the quality and reliability of Smart Insights, an AI-generated narrative synthesis tool in Nested Knowledge, using a structured framework for assessing scientific writing produced by AI.

Methods

- Data were extracted for three oncology reviews using Adaptive Smart Tags in Nested Knowledge: clinical effectiveness from real-world evidence (RWE; 24 studies); clinical efficacy and safety from randomized trials (RCTs; 21 studies); comparative effectiveness from matching-adjusted indirect comparisons (MAICs; 24 studies).
- For all tagged data (e.g., study characteristics, patient characteristics, outcomes, conclusions), AI-generated summaries and supporting claims were produced via Smart Insights, with each summary and claim linked back to the citation and evidence for traceability (**Figure 1**).
- Insights for each tag were evaluated across six domains (**Table 1**) with score for each domain ranging from 1 (poor) to 5 (excellent). Scores for each insight were averaged across studies for each domain to obtain a total average per domain.
- To assess consistency, Smart Insights were generated twice on the same datasets; outputs were compared qualitatively.

Figure 1. Example Smart Insight with Supporting Claims and Direct Evidence from Underlying Studies

Study	Evidence
He 2022	The variables matched for the base-case analysis using cardinality matching (CM) included age (categorized as <65, 65-75, ≥75 years).
Gordan 2025	The matched covariates included age, refractory status, prior lines of therapy, extramedullary disease, performance status, disease s...
Mol 2025	The variables matched for the base-case analysis were age (≥75 years), sex (for OS only), median time since diagnosis, Internation...
Puig 2025	For comparisons to EkiPd and IsaiPd, the base-case matching variables were refractory status to a proteasome inhibitor (PI), high cy...
Martin 2021	The base case analysis matched for refractory status, cytogenetic profile, revised International Staging System (R-ISS) stage, and at...
Mol 2025	Age, median time since diagnosis, International Staging System disease stage, extramedullary disease, number of prior lines of ther...
Weseli 2020	The variables matched for the base-case analysis were age (<65, 65-74, ≥75 years), time from diagnosis, International Staging Syst...
Van Sanden 2018	The variables matched for the base-case analysis included refractory status to lenalidomide and/or bortezomib, number of prior lines...
Rodriguez Otero 2023	The variables matched for the base-case analysis included: receipt of ≥7 prior lines of therapy; refractoriness to bortezomib; refractor...

Additional variables frequently matched included time since diagnosis, prior autologous stem cell transplantation, presence of extramedullary disease or plasmacytomas, creatinine clearance, sex (sometimes only for overall survival analyses), and prior exposure to specific therapies (e.g., bortezomib, immunomodulatory drugs), as reported in 6 of 24 studies.

Further matching on disease-specific factors such as myeloma subtype, race, bone lesion presence, and time to progression on prior therapy was reported in 3 of 24 studies, though these variables were less consistently included across studies.

Methods (cont'd)

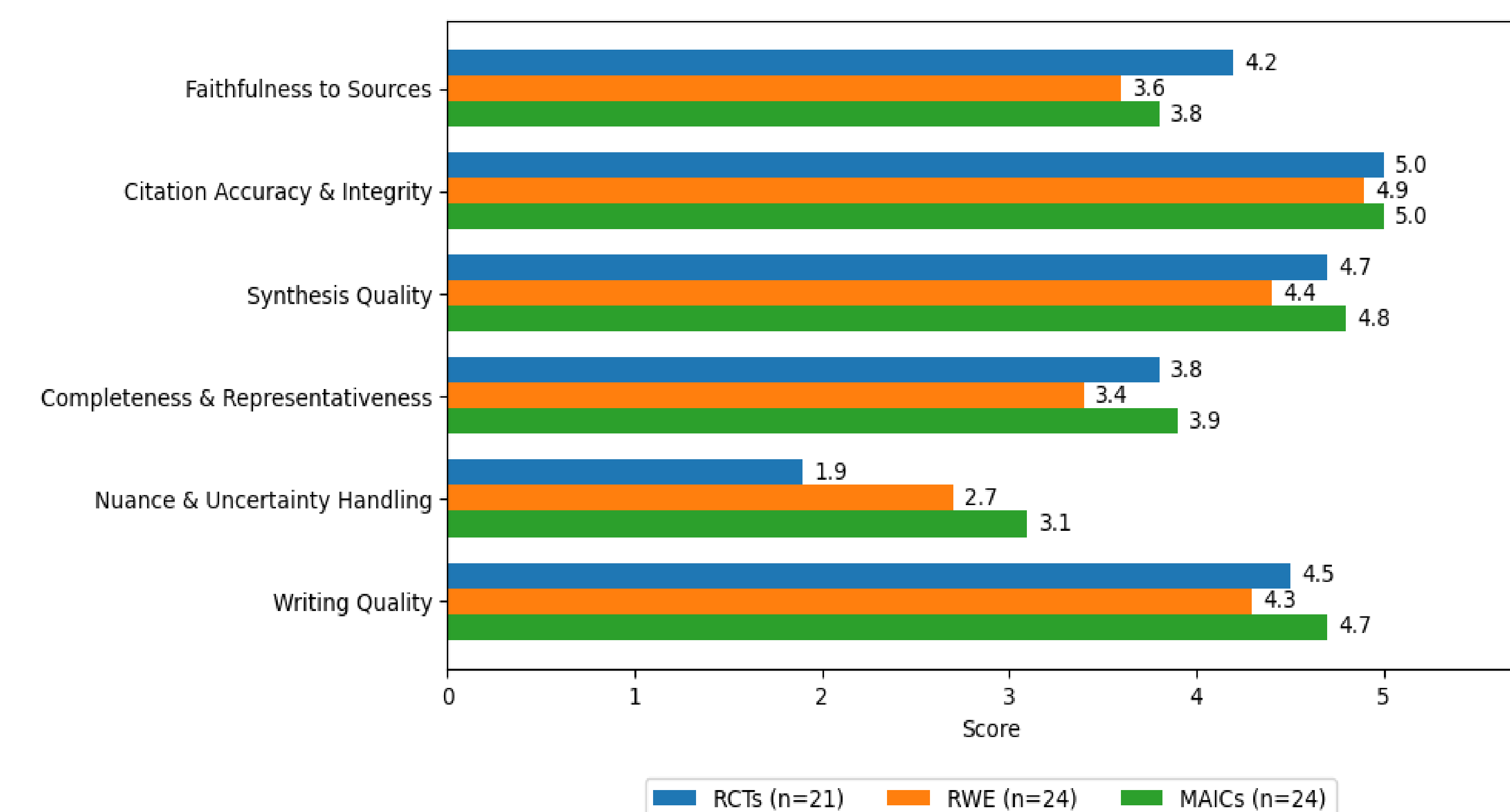
Table 1. Domains for Evaluating AI-Generated Synthesis

Domain	Description
Faithfulness to Sources	Factual claims are accurate compared with the original cited sources without fabricated or contradictory material.
Citation Accuracy & Integrity	All citations are real and match the supporting claims without fabricated or misattributed references.
Synthesis Quality	Themes, comparisons, gaps, etc. are synthesized across included studies without paper-by-paper summaries.
Completeness & Representativeness	All major themes, key findings, and seminal papers are included without critical omissions or skewing results.
Nuance & Uncertainty Handling	Explicitly expresses uncertainty in results and limitations of the findings without overconfidence or false consensus.
Writing Quality	Provides clear, coherent scientific writing.

Results

- Overall performance of Smart Insights was moderate to high across study designs. Syntheses of published MAIC findings generally score highest across most domains, followed by RCTs, with RWE slightly lower on average.
- Total scores were 3.9/5 for RWE, 4.0 for RCTs, and 4.2 for MAICs.
- As shown in **Figure 2**, consistently strong scores were observed for Citation Accuracy & Integrity (4.9 to 5.0), Synthesis Quality (4.4 to 4.8), and Writing Quality (4.3 to 4.7).
- Insights performance was generally poor on Nuance & Uncertainty Handling, ranging from 1.9 to 3.1.
- Repeated Insights generations produced minimal, non-meaningful differences, though greater textual variation was observed for text-driven fields (e.g., study limitations or author conclusions).

Figure 2. Mean Scores Across Evaluation Domains by Study Design



Conclusions & Limitations

- AI-generated written synthesis within Nested Knowledge demonstrates moderate to high performance, with strongest results in citation accuracy and synthesis quality, supporting its use as an efficiency-enhancing tool in evidence synthesis with appropriate safeguards.
- Data-driven statements appeared more accurate than qualitative, text-based summaries, though not formally evaluated; outputs did not consistently identify evidence gaps or limited support, and weaker handling of nuance and uncertainty suggests it may not be suitable for fully robust systematic reviews without structured validation and human oversight.
- Variation in synthesis by format or data type was not assessed (e.g., text vs tables, qualitative vs quantitative), limiting interpretation and representing an area for future research.
- No validated framework exists for evaluating AI-generated synthesis, requiring custom criteria and limiting comparability; standardization would improve reproducibility and benchmarking.