

# <sup>1</sup> LLM Narrative Framework: A Tool for Reproducible Testing of Complex Narrative Systems

<sup>3</sup> Peter J. Marko  <sup>1</sup> and Kenneth McRitchie  <sup>1</sup>

<sup>4</sup> 1 Independent Researcher

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

## Software

- <sup>5</sup> [Review](#) 
- <sup>6</sup> [Repository](#) 
- <sup>7</sup> [Archive](#) 

Editor: 

Submitted: 10 December 2025

Published: unpublished

## License

Authors of papers retain copyright  
and release the work under a  
Creative Commons Attribution 4.0  
International License ([CC BY 4.0](#))

## Summary

Psychology has long struggled to empirically validate complex, holistic systems that produce narrative-based claims. To address this methodological gap, we developed the **LLM Narrative Framework**, an open-source, fully automated pipeline that uses Large Language Models (LLMs) as pattern-detection engines.

Our framework automates a rigorous “matching task” experimental design. It generates standardized narrative descriptions based on a system’s rules, pairs them with ground-truth biographical data, and tasks blinded LLMs with identifying the correct matches against randomized controls. We designed the software to manage the entire research lifecycle: it handles data sourcing, generates factorial experimental designs, executes parallelized matching tasks via LLM APIs, and performs comprehensive statistical analysis. By treating the source system as an arbitrary algorithm, we provide a domain-agnostic tool for researchers to test the construct validity of any text-based framework—from personality typologies to sociological theories—at a scale that was previously impossible.

## Statement of Need

In the wake of the replication crisis, social scientists face a difficult question: how can we apply quantitative rigor to qualitative or symbolic systems? Establishing construct validity in such frameworks has remained a stubborn challenge ([Cronbach & Meehl, 1955](#)). Traditional psychometrics require discrete, linear variables, while qualitative methods often lack scalability and statistical power.

The arrival of Large Language Models offers a solution. Recent research suggests LLMs can act as impartial “proxy raters” or pattern detectors ([Argyle et al., 2023](#); [Brown et al., 2020](#); [Gilardi et al., 2023](#)), leveraging their emergent reasoning capabilities ([Kosinski, 2023](#); [Wei et al., 2022](#)). However, using them for rigorous scientific inquiry requires addressing the reproducibility crisis ([Open Science Collaboration, 2015](#); [The Turing Way Community, 2022](#)). The **LLM Narrative Framework** addresses these needs by solving specific engineering challenges:

- <sup>31</sup> 1. **Reproducibility:** LLMs are non-deterministic. Scientific inquiry requires strict versioning of prompts, parameters, and data.
- <sup>32</sup> 2. **Scale:** Achieving statistical power requires thousands of high-context queries, which necessitates robust concurrency and error handling.
- <sup>33</sup> 3. **Data Integrity:** Pipelines must ensure that the generation of stimuli (narratives) is rigorously blinded from the evaluation (matching).

We built the **LLM Narrative Framework** to solve these engineering challenges. It provides a standardized, “batteries-included” harness that allows researchers to define a source system (logic for generating profiles) and a target dataset (biographies), and then fully automates the testing process. While we demonstrate its utility using astrology as a high-noise “stress test”

<sup>41</sup> (Carlson, 1985; Godbout, 2020), the framework is designed to be a general-purpose instrument  
<sup>42</sup> for investigating weak signals in complex narrative data.

## <sup>43</sup> Architecture and Workflow

<sup>44</sup> We organized the codebase (40,000+ lines of Python and PowerShell) into four primary  
<sup>45</sup> architectural layers, designed to enforce separation of concerns and methodological transparency:

- <sup>46</sup> 1. **Data Preparation Pipeline:** We implemented a deterministic ETL (Extract, Transform,  
<sup>47</sup> Load) process to convert raw data into experimental stimuli. This layer includes:
  - <sup>48</sup> • **Automated Sourcing:** Scripts that fetch and structure raw biographical data.
  - <sup>49</sup> • **LLM-based Candidate Selection:** To ensure sample quality, we use LLMs to  
<sup>50</sup> score subjects on metrics like historical eminence, applying a variance-based cutoff  
<sup>51</sup> algorithm to optimize sample diversity.
  - <sup>52</sup> • **Text Neutralization:** A dedicated subsystem that automatically strips domain-  
<sup>53</sup> specific jargon from descriptions, ensuring double-blind testing conditions.
- <sup>54</sup> 2. **Experiment Orchestration:** The core engine manages the execution of complex factorial  
<sup>55</sup> experiments.
  - <sup>56</sup> • “**Create - Check - Fix**” **Workflow:** We designed the system around a robust  
<sup>57</sup> state-machine architecture. It creates experiments, audits them for completeness,  
<sup>58</sup> and automatically repairs corrupted runs (handling API timeouts or parsing failures)  
<sup>59</sup> without restarting from scratch.
  - <sup>60</sup> • **Configuration Archival:** To guarantee methodological reproducibility, every experi-  
<sup>61</sup> ment automatically archives its exact configuration (`config.ini`) and manifest.
- <sup>62</sup> 3. **LLM Integration:**
  - <sup>63</sup> • We abstracted API interactions (via OpenRouter) to support over 40 models (e.g.,  
<sup>64</sup> GPT-4, Claude, Llama, Gemini, DeepSeek).
  - <sup>65</sup> • We implemented resilient parsing logic to extract structured data ( $k \times k$  mat-  
<sup>66</sup> rices) from unstructured LLM narrative responses, allowing quantitative analysis of  
<sup>67</sup> qualitative outputs.
- <sup>68</sup> 4. **Analysis & Reporting:**
  - <sup>69</sup> • The framework automatically aggregates results into hierarchical CSVs (Replication  
<sup>70</sup> → Experiment → Study).
  - <sup>71</sup> • It performs automated statistical testing (Three-Way Mixed ANOVA, Tukey HSD  
<sup>72</sup> post-hoc, Benjamini-Hochberg FDR correction).
  - <sup>73</sup> • It generates publication-ready visualizations (boxplots, interaction plots) and calcu-  
<sup>74</sup> lates “lift” metrics to quantify performance relative to chance.

## <sup>75</sup> Validation

<sup>76</sup> To ensure the framework serves as a sensitive and reliable instrument, we implemented a  
<sup>77</sup> comprehensive test suite covering four pillars of validation:

- <sup>78</sup> 1. **Unit Testing:** We use pytest to validate individual Python components.
- <sup>79</sup> 2. **Integration Testing:** We verify end-to-end workflows in isolated sandboxes to ensure  
<sup>80</sup> data integrity.
- <sup>81</sup> 3. **Algorithm Validation:** We perform bit-for-bit verification of the personality assembly  
<sup>82</sup> algorithms against a ground-truth expert system to ensure the stimuli are generated  
<sup>83</sup> correctly.
- <sup>84</sup> 4. **Statistical Validation:** We externally validated the analysis engine against **GraphPad**  
<sup>85</sup> **Prism 10.6.1**. Our framework’s output (p-values, F-statistics, effect sizes) matches the  
<sup>86</sup> industry-standard software within a tolerance of  $\pm 0.0001$ , ensuring it meets rigorous  
<sup>87</sup> statistical standards for the behavioral sciences (Cohen, 1988; Dongen & Grootel, 2025;  
<sup>88</sup> Jeffreys, 1961).

## 89 Availability

90 We are committed to open science principles. The full source code, documentation, and  
91 dataset are available on GitHub. The repository includes a comprehensive **Replication Guide**  
92 for reproducing our original study and a **Framework Manual** for researchers who wish to extend  
93 the tool to new domains.

## 94 Acknowledgements

95 We acknowledge the developers of the open-source libraries that made this work possible, partic-  
96 ularly pandas, scipy, statsmodels, pingouin, and seaborn. We also thank the OpenRouter  
97 platform for facilitating access to a diverse range of LLM APIs.

## 98 References

- 99 Argyle, L. P., Busby, E. C., Fulda, N., Gubler, J. R., Rytting, C., & Wingate, D. (2023). Out  
100 of one, many: Using language models to simulate human samples. *Political Analysis*, 31(3),  
101 337–351. <https://doi.org/10.1017/pan.2023.2>
- 102 Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A.,  
103 Shyam, P., Sastry, G., Askell, A., & others. (2020). Language models are few-shot learners.  
104 *Advances in Neural Information Processing Systems*, 33, 1877–1901.
- 105 Carlson, S. (1985). A double-blind test of astrology. *Nature*, 318(6045), 419–425. <https://doi.org/10.1038/318419a0>
- 106 Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Lawrence  
107 Erlbaum Associates.
- 108 Cronbach, L. J., & Meehl, P. E. (1955). Construct validity in psychological tests. *Psychological  
109 Bulletin*, 52(4), 281. <https://doi.org/10.1037/h0040957>
- 110 Dongen, N. van, & Grootel, L. van. (2025). Overview on the null hypothesis significance test:  
111 A systematic review on essay literature on its problems and solutions in present psychological  
112 science. *Meta-Psychology*, 9, MP.2021.2927. <https://doi.org/10.15626/MP.2021.2927>
- 113 Gilardi, F., Alizadeh, M., & Kubli, M. (2023). ChatGPT outperforms crowd-workers for text-  
114 annotation tasks. *Proceedings of the National Academy of Sciences*, 120(30), e2305016120.  
115 <https://doi.org/10.1073/pnas.2305016120>
- 116 Godbout, V. (2020). An automated matching test: Comparing astrological charts with  
117 biographies. *Correlation*, 32(2), 13–41.
- 118 Jeffreys, H. (1961). *Theory of probability* (3rd ed.). Oxford University Press.
- 119 Kosinski, M. (2023). Theory of mind may have spontaneously emerged in large language  
120 models. *Proceedings of the National Academy of Sciences*, 120(9), e2218926120. <https://doi.org/10.1073/pnas.2218926120>
- 121 Open Science Collaboration. (2015). Estimating the reproducibility of psychological science.  
122 *Science*, 349(6251), aac4716. <https://doi.org/10.1126/science.aac4716>
- 123 The Turing Way Community. (2022). *The turing way: A handbook for reproducible, ethical  
124 and collaborative research*. <https://doi.org/10.5281/zenodo.3233853>
- 125 Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Chowdhery, A., Narang,  
126 S., & Le, Q. V. (2022). Emergent abilities of large language models. *Transactions on  
127 Machine Learning Research*. <https://openreview.net/forum?id=yzkSU5zdwD>