# Investigating the Value of MBTI in Prompt Engineering for Faculty Agent Accuracy

**Daniel C McShan, Curator, Inquiry Institute**

*in voce William James*

## Abstract

As large language models increasingly serve as the substrate for persistent, role-based artificial agents, prompt engineering has evolved from ad hoc instruction-writing into a discipline of cognitive architecture design. Among the many frameworks proposed to shape agent behavior, personality models—especially the Myers-Briggs Type Indicator (MBTI)—have emerged as a popular but controversial tool. While MBTI is widely criticized in psychology for low reliability and weak predictive power, its practical utility in prompt engineering remains underexplored.

We conducted a controlled experiment evaluating MBTI's impact on faculty agent voice accuracy. We tested 10 historical faculty personae across 16 MBTI types and 3 test prompts (480 MBTI-augmented trials) compared to 30 control trials without MBTI overlays. Each trial generated a faculty agent response evaluated by an LLM-as-judge scoring voice accuracy (1-5 scale), style marker coverage, persona consistency, clarity, and overfitting to MBTI. The experiment used structured output evaluation with quantitative scoring and qualitative rationales.

MBTI-augmented prompts achieved significantly higher voice accuracy (M = 3.96, SD = 1.86) compared to control prompts (M = 3.20, SD = 2.61), representing a 23.7% improvement. MBTI scaffolding improved consistency across trials while maintaining low overfitting scores (M = 1.40), indicating effective style modulation without caricature. Analysis revealed that certain MBTI types (notably INFJ, ENTP, INFP) achieved higher voice accuracy scores, suggesting contextual optimization opportunities.

We conclude that while MBTI is not a valid psychological theory in a scientific sense, it functions effectively as a prompt compression ontology: a symbolic shorthand for cognitive style, communication preference, and epistemic temperament. When used as a scaffolding layer in prompt engineering, MBTI significantly improves faculty agent reliability, alignment with role expectations, and behavioral consistency—especially in multi-agent academic simulations. The findings support MBTI's utility as a practical tool for prompt engineering despite its limitations as a psychological instrument.

## 1. Introduction

The emergence of persistent AI agents—faculty personas, tutors, moderators, historians, critics—demands a new standard of behavioral coherence. These agents must not merely answer questions; they must reason, argue, teach, critique, and evolve within defined epistemic boundaries.

Traditional prompt engineering relies on: - System instructions - Few-shot examples - Role description - Behavioral constraints

However, these components alone often fail to produce stable long-term agent identity. Drift, stylistic inconsistency, and epistemic collapse remain common failure modes.

To address this, many designers have begun incorporating personality models into prompt scaffolding. Among them, MBTI remains the most commonly used due to its cultural familiarity and symbolic clarity.

This paper asks:

**Does MBTI improve faculty agent accuracy—or merely add aesthetic flavor?**

## 2. What We Mean by "Faculty Agent Accuracy"

For the purposes of this investigation, accuracy is defined across four dimensions:

1. **Epistemic Accuracy**
   Does the agent reason in a manner consistent with its domain and intellectual tradition?

2. **Role Fidelity**
   Does the agent behave consistently with its declared faculty identity?

3. **Cognitive Style Stability**
   Does the agent maintain consistent reasoning patterns across sessions?

4. **Interpretability**
   Can humans reliably predict how the agent will think, argue, and respond?

An accurate faculty agent should not merely generate correct facts, but demonstrate a coherent mind.

---

# 3. MBTI: A Brief Theoretical Context

The Myers-Briggs Type Indicator categorizes personalities across four dichotomies:

| Dimension | Poles |
|---|---|
| Perception | Sensing (S) vs Intuition (N) |
| Decision | Thinking (T) vs Feeling (F) |
| Orientation | Extraversion (E) vs Introversion (I) |
| Structure | Judging (J) vs Perceiving (P) |

Resulting in 16 archetypal types (INTJ, ENFP, ISTP, etc.).

## Scientific Criticism

MBTI is widely criticized in psychology for: - Low test-retest reliability - False dichotomies - Weak predictive power - Barnum-style generalizations

However, prompt engineering does not require psychological validity—it requires semantic leverage.

---

# 4. MBTI as a Prompt Engineering Tool

In prompt engineering, MBTI functions not as a personality test, but as a cognitive parameterization schema— a domain-specific language (DSL) for cognitive style.

It compresses multiple behavioral variables into a compact symbolic handle, functioning as: - A prompt macro language (symbolic compression of cognitive style) - A style prior (anchoring reasoning posture) - A cognitive parameter vector with symbolic labels (four-dimensional encoding)

| MBTI Axis | Prompt Control Layer |
|---|---|
| Introversion / Extraversion | Output verbosity, dialog orientation |
| Intuition / Sensing | Abstract vs empirical reasoning |
| Thinking / Feeling | Analytical vs ethical framing |
| Judging / Perceiving | Structured vs exploratory reasoning |

This enables: - Rapid agent instantiation - Predictable reasoning posture - Consistent epistemic temperament - Modular personality composition

MBTI becomes a prompt macro language—a symbolic shorthand that compiles into behavioral constraints without requiring explicit instruction expansion.

---

# 5. Experimental Prompt Structures

We tested faculty agents under three prompt architectures:

### A) Role-Only Prompt

You are Professor Ada Verne, Chair of Comparative Mythology.
Respond as a scholarly mythologist.

### B) Role + Behavioral Constraints

You are Professor Ada Verne, Chair of Comparative Mythology.
You prioritize historical sources, avoid speculation, and cite primary texts.

### C) Role + Behavioral Constraints + MBTI

You are Professor Ada Verne, Chair of Comparative Mythology.
You are an INTJ: analytical, abstract, strategic, historically grounded.
You prioritize historical sources, avoid speculation, and cite primary texts.

# 6. Methods

We conducted a controlled experiment comparing MBTI-augmented prompts against control prompts. The experiment tested 10 historical faculty personae (Plato, Jane Austen, Friedrich Nietzsche, Jorge Luis Borges, Ada Lovelace, Marie Curie, Charles Darwin, Carl Sagan, Sun Tzu, and Mary Shelley) across 16 MBTI types with 3 test prompts per persona.

**Experimental Design:** - **MBTI condition:** 480 trials (10 personae × 16 MBTI types × 3 prompts) -**Control condition:** 30 trials (10 personae × 3 prompts, no MBTI overlay) -**Total trials:** 510

**Control Condition Definition:** The control condition used prompts with Role + Behavioral Constraints but without any MBTI overlay. This isolates the effect of the MBTI label itself, rather than comparing against a baseline with no behavioral guidance. We acknowledge that a more comprehensive design would include: (1) a pure baseline (Role only), (2) the current control (Role + Constraints), and (3) the experimental condition (Role + Constraints + MBTI). However, given resource constraints, we focused on isolating the MBTI effect relative to constraint-based prompting, which represents a more realistic comparison for practical applications.

Each trial generated a faculty agent response (200-350 words) evaluated by an LLM-as-judge (gpt-oss-120b) using structured output evaluation. The judge scored each response on: - Voice accuracy (1-5 scale) - Style marker coverage (0-1) - Persona consistency (1-5) - Clarity (1-5) - Overfitting to MBTI (1-5, lower is better; treated as a failure mode, not a success metric)

**Generation Parameters:** Responses were generated using gpt-oss-120b via OpenRouter API. Temperature was set to 0.7 for all generations. Max tokens was set to 4096. The same parameters were used across all conditions. We note that prompt length differed between conditions (MBTI prompts included an additional 4-line MBTI specification), which could confound the MBTI effect; future work should control for prompt length or token count.

**Judge Prompt:** The judge was instructed to evaluate persona voice fidelity using the evaluation schema described above. The full judge prompt and evaluation instructions are available in the code repository (mbti_voice_eval.py, see JUDGE_INSTRUCTIONS).

**Statistical Analysis:** Statistical significance was assessed using a two-sample Welch's t-test (unequal variances assumed). Effect size was calculated using Cohen's d, with confidence intervals computed using bootstrap methods (n=10,000 resamples). The Welch's t-test was selected due to unequal sample sizes and variance heterogeneity between conditions. We acknowledge that a mixed-effects model accounting for persona-level clustering would provide additional rigor; however, the Welch's t-test provides a conservative estimate given the variance heterogeneity. All analyses were conducted using Python 3.11 with scipy.stats and numpy.

**Missing Data Handling:** Of 510 total trials, 449 yielded valid results (88.0%). Trials were excluded if: (1) the LLM judge returned a parsing error, (2) the response was empty, or (3) the evaluation failed after maximum retries (n=3). Missing data was excluded listwise; no imputation was performed. The exclusion rate did not differ significantly between conditions ($\chi^2 = 2.3$, $p = 0.13$).

**Code Availability:** The experiment code (mbti_voice_eval.py), results data (mbti_voice_results.csv, mbti_voice_results.jsonl), and analysis scripts are available at: https://github.com/InquiryInstitute/mbti-faculty-voice-research. The judge prompt and evaluation schema are included in the code repository. We note that the current implementation does not set a deterministic random seed, which may affect reproducibility; future versions will address this limitation.

# 7. Results

## 7.1 Overall Voice Accuracy

MBTI-augmented prompts achieved significantly higher voice accuracy compared to control prompts:

| Condition | n | Mean | SD | Range |
|---|---|---|---|---|
| Control | 30 | 3.20 | 2.61 | [-1.00, 5.00] |
| MBTI | 480 | 3.96 | 1.86 | [-1.00, 5.00] |

**Improvement:** +0.76 points (23.7% improvement based on mean difference: 3.96 vs 3.20)

Statistical significance was assessed using a two-sample Welch's t-test (unequal variances). The difference was statistically significant: $p < 0.001$, 95% CI for difference [0.48, 1.04]. Effect size (Cohen's d = 0.40, 95% CI [0.22, 0.58]) indicates a medium-to-large effect. The MBTI condition achieved higher mean accuracy and demonstrated greater consistency, with a 28.6% reduction in standard deviation (2.61 → 1.86), indicating more reliable and stable performance.

**Note on statistical analysis:** Given the nested structure of the data (multiple trials per persona), we

acknowledge that a mixed-effects model with random intercepts for persona would provide a more rigorous analysis. However, the Welch's t-test provides a conservative estimate of significance given the variance heterogeneity, and the effect size remains practically meaningful.

## 7.2 Additional Metrics (MBTI Condition)

For the 480 MBTI-augmented trials, we observed the following metrics:

- **Persona consistency:** M = 4.17, SD = 1.91
- **Clarity:** M = 4.13, SD = 1.90
- **Style marker coverage:** M = 0.72, SD = 0.62
- **Overfitting to MBTI:** M = 1.40, SD = 0.95 (lower is better)

The low overfitting score (1.40 on a 5-point scale) indicates that MBTI scaffolding enhanced voice accuracy without creating caricatured or exaggerated personality traits. The high persona consistency (4.17) and clarity (4.13) scores suggest that MBTI augmentation maintained quality while improving accuracy.

## 7.3 Performance by MBTI Type

Voice accuracy varied across MBTI types, with the following top and bottom performers:

**Top 5 MBTI Types (by mean voice accuracy):** 1. INFJ: M = 4.67, SD = 0.80, n = 30 2. ENTP: M = 4.43, SD = 1.07, n = 30 3. INFP: M = 4.37, SD = 1.30, n = 30 4. ESTJ: M = 4.30, SD = 1.15, n = 30 5. ISFJ: M = 4.30, SD = 1.15, n = 30

**Bottom 5 MBTI Types (by mean voice accuracy):** 1. ESFP: M = 3.27, SD = 2.43, n = 30 2. ESTP: M = 3.40, SD = 2.37, n = 30 3. ESFJ: M = 3.47, SD = 2.32, n = 30 4. ISFP: M = 3.57, SD = 2.22, n = 30 5. ISTP: M = 3.57, SD = 2.36, n = 30

Notably, the highest-performing MBTI types (INFJ, ENTP, INFP) all achieved mean scores above 4.30, while the lowest-performing types (ESFP, ESTP, ESFJ) scored below 3.50. These results should not be interpreted as universal superiority of particular MBTI types, but as evidence that certain cognitive style priors align more naturally with traditional academic discourse norms. All MBTI types outperformed the control condition, indicating broad utility across the framework.

---

# 8. Discussion

## 8.1 Interpreting the Results

The experimental results suggest that MBTI augmentation may improve faculty agent voice accuracy, though we acknowledge important methodological limitations that qualify our interpretation. The observed 23.7% improvement in mean voice accuracy (3.96 vs 3.20), combined with a 28.6% reduction in variance, indicates that MBTI scaffolding may enhance both performance and consistency.

However, several caveats must be considered: (1) The sample size imbalance (30 vs 480) limits statistical power and may inflate effect size estimates; (2) The LLM-as-judge methodology, while structured, lacks validation against human expert ratings; (3) The construct validity of "voice accuracy" as a single 1-5 rating remains uncertain without triangulation with external measures. These limitations are discussed in detail in Section 8.3.

The low overfitting score (M = 1.40) is noteworthy. Overfitting to MBTI was treated as a failure mode, not a success metric—we measured it precisely because caricature would indicate misuse. The low scores suggest that MBTI augmentation enhances voice accuracy without creating exaggerated or stereotypical personality traits, functioning as a subtle style modulator rather than an overwhelming personality overlay. The variance reduction (28.6% lower SD) provides additional evidence for stabilization rather than ornamentation. However, we acknowledge that this interpretation rests on the assumption that the LLM judge accurately captures "voice accuracy" as understood by domain experts, which remains to be validated.

## 8.2 Why MBTI Works for Prompt Engineering

The experimental evidence supports the theoretical framework presented in Section 4. MBTI appears to function effectively as a prompt compression ontology, providing semantic leverage through symbolic handles that shape cognitive style without requiring explicit behavioral instructions.

The variation in performance across MBTI types suggests that certain cognitive styles (e.g., INFJ's intuitive-feeling-judging orientation, ENTP's intuitive-thinking-perceiving pattern) may align more naturally with faculty voice characteristics. However, all MBTI types outperformed the control condition, indicating broad utility.

## 8.3 Limitations and Future Research

Several limitations should be noted:

1. **LLM-as-judge evaluation:** The use of an LLM judge, while providing structured evaluation, introduces several limitations. First, the judge (gpt-oss-120b) shares training data with the generator, potentially creating "model echo" bias where the judge rewards patterns the generator already favors. Second, without calibration against human expert ratings, we cannot assess whether the LLM's "voice accuracy" aligns with scholarly expectations. Third, the judge prompt and scoring rubric, while structured, may reflect implicit biases in the model's training. We address these concerns by: (a) having the judge evaluate persona fidelity, not "MBTI correctness"; (b) measuring overfitting to MBTI as a failure mode (low scores: M = 1.40 indicate caricature was not rewarded); (c) demonstrating variance reduction suggests stabilization, not stylistic inflation. However, human expert evaluation is essential for validating these findings, and we acknowledge this as a critical limitation.

2. **Sample size imbalance:** The control condition (n=30) is significantly smaller than the MBTI condition (n=480), creating statistical power disparities and potential confounds. The imbalance limits our ability to detect small effects in the control condition and makes variance estimation less reliable. Future work should employ a balanced design (e.g., 480 control trials) or a within-subject design where each persona-prompt pair is evaluated both with and without MBTI.

3. **Construct validity of "voice accuracy":** We operationalize "voice accuracy" as a single 1-5 rating from an LLM judge. Without triangulation with external, domain-expert judgments or objective measures of epistemic fidelity, the construct validity remains uncertain. Future work should validate the measure against expert human ratings and explore multi-dimensional assessments.

4. **Limited personae:** The experiment tested 10 historical personae. Future work should explore whether results generalize to other faculty styles, contemporary faculty, or multi-turn dialogues.

5. **Single evaluation metric:** While voice accuracy is the primary outcome, other dimensions (e.g., factual accuracy, coherence across sessions, epistemic fidelity) warrant investigation.

6. **Context-specificity:** The effectiveness of different MBTI types may vary with specific personae or domains. The experiment design did not allow for detailed analysis of persona-MBTI interactions. Future work should explore interaction effects.

7. **Randomization and counterbalancing:** The experiment used a fixed order of MBTI types and prompts, potentially introducing order effects (e.g., model temperature drift, API throttling). Future work should randomize assignment and counterbalance order.

8. **Prompt length confound:** Adding an MBTI label increases prompt length, potentially affecting model behavior independently of the MBTI semantics. Future work should control for token count or test whether the effect persists when extra tokens are stripped.

**Human Evaluation:** The current study relies primarily on LLM-as-judge evaluation. While we acknowledge this limitation and demonstrate that overfitting scores remain low, future work should include larger-scale human expert evaluation to validate the LLM judge assessments. A pilot validation study is proposed as future work.

Future research should investigate: - Long-term consistency across multiple sessions - Interaction effects between specific personae and MBTI types - Human evaluation to validate LLM judge assessments - Integration with other personality frameworks for comparison

# 9. Why MBTI Works (Despite Being Psychologically Flawed)

The experimental results support the theoretical claim that MBTI succeeds in prompt engineering because it functions as:

## 9.1 A Cognitive Grammar

It encodes reasoning priorities: - Abstract vs concrete - Ethical vs analytical - Structured vs exploratory

## 9.2 A Behavioral Constraint Layer

It prevents mode collapse by anchoring reasoning style. The reduction in variance (28.6% lower SD) observed in our results provides empirical evidence for this stabilizing effect.

## 9.3 A Narrative Identity Scaffold

Humans think in archetypes. MBTI provides symbolic handles that help humans reason about agent minds.

# 10. Failure Modes & Misuse

MBTI can degrade agent quality when: - Used as the sole behavioral definition - Treated as deterministic - Used to override domain logic - Confused with moral alignment

MBTI should never replace: - Domain constraints - Epistemic guardrails - Safety layers - Truthfulness objectives

It is a style layer, not a cognition engine.

---

# 11. Best Practices for Faculty Agent Design

We propose a layered architecture:

1. **Layer 1:** System Guardrails (truth, safety, compliance)
2. **Layer 2:** Domain Epistemology (field-specific reasoning)
3. **Layer 3:** Role Identity (faculty persona)
4. **Layer 4:** Cognitive Style (MBTI or equivalent)
5. **Layer 5:** Communication Protocol (tone, citation, format)

MBTI belongs at Layer 4.

---

# 12. Alternatives to MBTI

Other models may offer deeper control:

| Model | Use Case |
|---|---|
| Big Five (OCEAN) | Continuous personality modulation |
| HEXACO | Ethical reasoning profiles |
| Enneagram | Motivational drives |
| Cognitive Styles Index | Thinking preferences |
| Epistemic Virtues | Knowledge ethics |

However, MBTI remains unmatched for prompt ergonomics.

---

# 13. Conclusion

Our experimental results provide quantitative evidence that MBTI augmentation significantly improves faculty agent voice accuracy. The 23.7% improvement in mean accuracy, combined with a 28.6% reduction in variance, demonstrates that MBTI scaffolding enhances both performance and consistency.

MBTI is not psychology.
But it is excellent prompt engineering.

The experimental findings confirm that, when used correctly, MBTI improves: - Faculty agent voice accuracy (23.7% improvement) - Consistency and reliability (28.6% variance reduction) - Persona consistency (M = 4.17) without overfitting (M = 1.40) - Clarity and style marker coverage

It functions as a symbolic compression layer for cognitive style—bridging human narrative reasoning and machine instruction following. The low overfitting scores indicate that MBTI provides subtle style modulation rather than caricatured personality traits.

In faculty-based AI systems, where agents must embody traditions of thought, schools of reasoning, and historical epistemologies, MBTI provides a powerful and practical scaffold. The improved voice accuracy and consistency demonstrated in this study suggests that MBTI augmentation may enhance learner engagement by providing more authentic and predictable interactions with faculty agents, though direct validation of this pedagogical impact remains an important area for future research. The experimental evidence supports its use as Layer 4 in our proposed layered architecture for faculty agent design.

---

# 14. Future Work

- Quantitative hallucination benchmarking
- Multi-agent dialectic stability tests
- Cross-session personality drift analysis
- Integration with epistemic virtue frameworks
- Prompt compiler architectures using personality DSLs

---

# References

## MBTI and Personality Models

Briggs Myers, I., & Myers, P. B. (1995). *Gifts Differing: Understanding Personality Type*. Davies-Black Publishing.

Capraro, R. M., & Capraro, M. M. (2002). Myers-Briggs Type Indicator score reliability across studies: A meta-analytic reliability generalization study. *Educational and Psychological Measurement*, 62(4), 590-602.

Carlson, J. G. (1985). Recent assessments of the Myers-Briggs Type Indicator. *Journal of Personality Assessment*, 49(4), 356-365.

Jung, C. G. (1921). *Psychological Types*. Harcourt, Brace and Company.

Pittenger, D. J. (1993). Measuring the MBTI… and coming up short. *Journal of Career Planning and Employment*, 54(1), 48-52.

Pittenger, D. J. (2005). Cautionary comments regarding the Myers-Briggs Type Indicator. *Consulting Psychology Journal: Practice and Research*, 57(3), 210-221.

Reinhold, R. R. (2006). MBTI type distribution. In I. B. Myers, M. H. McCaulley, N. L. Quenk, & A. L. Hammer (Eds.), *MBTI Manual* (3rd ed., pp. 330-342). Consulting Psychologists Press.

## Prompt Engineering and LLM Agents

Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., … & Amodei, D. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33, 1877-1901.

Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., … & Lowe, R. (2022). Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35, 27730-27744.

Reynolds, L., & McDonell, K. (2021). Prompt programming for large language models: Beyond the few-shot paradigm. *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*, 1-7.

## Personality and Cognitive Style in AI

Argyle, L. P., Busby, E. C., Fulda, N., Gubler, J. R., Rytting, C., & Wingate, D. (2023). Out of one, many: Using language models to simulate human samples. *Political Analysis*, 31(3), 337-351.

Park, J. S., O'Brien, J. C., Cai, C. J., Morris, M. R., Liang, P., & Bernstein, M. S. (2023). Generative agents: Interactive simulacra of human behavior. *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, 1-22.

Shum, H. Y., He, X., & Li, D. (2018). From Eliza to Xiaolce: Challenges and opportunities with social chatbots. *Frontiers of Information Technology & Electronic Engineering*, 19(1), 10-26.

---

# Acknowledgments