

What AI Found, What We Already Knew

Chungil Chae

Fri, 3 April 2026

What AI Found, What We Already Knew

Interpreting a Paper Autonomously Generated from Data by an AI-Scientist

Date written: 2026-04-03

Opening: We Handed the Data to a Machine

We ran 36 models. Across 5 methodological paradigms, drawing on 16 theories, we asked the same question of 282 cases three times over. We organized the results into two reports (a strategy edition and an interpretation edition).

Then we gave the same data to an **AI-Scientist**. “Analyze it yourself, and write a paper.”

The AI chose K-means clustering on its own. It autonomously carried out 5 experiments. It wrote an 8-page paper. The cost was \$8, the time 20 minutes.

What that AI found **overlaps exactly** with one piece of the conclusion we had already reached.

17 Chapter 1. What the AI Saw

18 The Choice of Clustering

19 The AI had no theoretical guidance whatsoever. It knew nothing of TOE, of P-E Fit, of absorptive
20 capacity. All we gave it were two seed_ideas (nonlinear extension of TOE, dose-response separation)
21 as hints.

22 The AI generated a third idea on its own: “**Classify firms with K-means, and analyze each**
23 **cluster separately.**” Its scores were Interestingness 9, Feasibility 9, Novelty 8. Among the four
24 ideas, it gave this one the highest rating of all—to itself.

25 This is the same intuition we had pursued with T2 (LPA) in our first-round analysis. The sense
26 that there are **hidden types** within the data. The suspicion that beneath the average lies
27 heterogeneity.

28 Three Profiles

29 The 3 clusters the AI classified:

Cluster	n	DT awareness	Smart systems	Firm size (log)	Q3_difference
0 (low readiness)	122	1.84	16.1	2.51	1.525
1 (high readiness)	51	3.24	25.4	3.34	1.498
2 (medium readiness)	70	1.89	21.3	3.73	1.247

30 Here the first odd point appears. **The Q3_difference of the high-readiness firms (cluster**
31 **1) is almost identical to that of the low-readiness firms (cluster 0).** 1.498 vs 1.525. DT
32 awareness is twice as high and smart systems are 60% greater, yet the training effect is the same.

33 This is exactly what we found in T2 (LPA) of the first-round analysis: “*A more prepared firm*
34 *does not necessarily learn better.*” It is the very pattern the interpretation report explained as the
35 “large-firm paradox” and the “ceiling effect.”

36 The Dramatic Gap of R^2 .002 \rightarrow .336

37 The AI’s most impactful finding:

- 38 • **Cluster 0** (low readiness, n=122): $R^2 = .002$. No variable explains the training effect.
- 39 • **Cluster 1** (high readiness, n=51): $R^2 = .336$. Smart systems (+) and firm size (-) are strong
40 predictors.

- 41 • **Cluster 2** (medium readiness, n=70): $R^2 = .018$. Again, nothing explains anything.
 42 What was $R^2 = .029$ across the entire dataset jumps to $.336$ once only the high-readiness firms are
 43 isolated. **12 times higher**. The AI called this “dramatic heterogeneity.”

44

45 Chapter 2. What We Already Knew

46 When the AI’s findings are set against our 36 analyses, the puzzle pieces fall into place.

47 Comparison 1: “Predictors Work Only in High Readiness”

48 The AI’s finding: R^2 is meaningful only in cluster 1.

49 Our findings: - **T2 (LPA)**: The difference in training effect across profiles was not significant, but
 50 the **mechanism** differed by profile. - **T3/T20 (QCA)**: Equifinality. There are multiple paths to
 51 a high training effect. \sim DT_AWARE (low awareness) is also a valid path. - **combo1 (RSM)**:
 52 The P-E Fit interaction term ($p=.002$) is significant only when “awareness and infrastructure are
 53 **both** high at the same time.”

54 What the AI found with K-means is **another expression** of what we had reached through QCA’s
 55 equifinality and RSM’s interaction term. “It works only in high readiness” and “the awareness \times
 56 infrastructure interaction term is significant” are different angles on the same phenomenon.

57 Comparison 2: “The Negative Effect of Firm Size”

58 The AI’s finding: in cluster 1, firm size = -0.839 ($p<.001$). Even among high-readiness firms, the
 59 larger the size, the **lower** the training effect.

60 Our findings: - **T14 (ANOVA)**: The large-firm paradox. An adverse effect in firms of large size. -
 61 **T15**: SF + multiple-participation reverse synergy. The limits of additional input in firms that
 62 already have a lot. - **RSM (B1)**: An inverted-U for firm size ($p=.029$), optimal at ~ 13 people.

63 The AI’s = -0.839 corresponds to the **right-hand declining segment of the inverted-U**.
 64 Because high-readiness (cluster 1) firms have a large average size ($\log 3.34 \sim 28$ people), they have
 65 already passed the optimum of the inverted-U (~ 13 people). They are in the region of diminishing
 66 returns.

67 This is **one side** of what the interpretation report explained as “the dual mechanism of Liability
 68 of Smallness + ceiling effect”—and the AI captured it.

69 **Comparison 3: “A Difference in Slopes, Not Intercepts”**

70 The AI’s most refined finding: adding cluster dummy variables yields non-significance ($p=.856$),
71 but adding cluster \times variable interaction terms yields significance ($p=.043$).

72 This means that while the **mean training effect of the three groups is similar** (no difference
73 in intercepts), **the mechanism that determines the training effect differs** (a difference in
74 slopes).

75 The integrative proposition of the interpretation report was precisely this:

76 “The effect of SME digital transformation training is not a matter of ‘present or absent,’
77 but appears ‘when awareness and infrastructure meet.’”

78 Not “present or absent” (difference in intercepts) \rightarrow consistent with the AI’s non-significant dummy
79 ($p=.856$). “When they meet” (conditional) \rightarrow consistent with the AI’s significant interaction term
80 ($p=.043$).

81

82 **Chapter 3. What the AI Missed**

83 The reason the AI’s paper received a Reject (an internal review—it reviews what it itself wrote,
84 based on a model built from CS-side papers) is not simply a journal mismatch. **The things**
85 **essential to a social-science paper are missing.**

86 **Missing 1: Theory**

87 The AI’s paper has no theory. There is a **description**—“there is heterogeneity”—but no **expla-**
88 **nation** of “why heterogeneity arises.”

89 We drew on 16 theories. P-E Fit explained the interaction term, diminishing returns the inverted-U,
90 equifinality the multiple paths, experiential learning the dose-response. **A finding without**
91 **theory is no more than a pattern.** That the AI found $R^2 .002 \rightarrow .336$ is impressive, but
92 without a theoretical answer to “why do predictors work only in high readiness?” it remains an
93 **unfinished finding.**

94 The answer: it is P-E Fit. Only in firms where the fit between awareness (P) and infrastructure
95 (E) is high does the mechanism activate by which additional organizational characteristics (smart
96 systems, firm size) predict the training effect. In firms with low fit, the mechanism itself is “switched
97 off.”

98 **Missing 2: Triangulation**

99 The AI reached its conclusion with K-means alone. We applied 5 methodologies to the same
100 question: - **LPA** (person-centered) → profile types - **QCA** (set-theoretic) → condition combinations
101 - **RSM** (nonlinear) → interaction terms / inverted-U - **DID/PSM** (causal inference) → treatment
102 effects - **CLPM** (longitudinal) → causal direction

103 A single K-means silhouette of .293 may be called “reasonable,” but **is that pattern confirmed**
104 **by other methodologies too?** Our answer was “yes, it converges across all 5.” The AI’s answer
105 remained at the level of “I’m not sure—I changed the seed and it came out similar.”

106 **Missing 3: The “Limits of Measurement” Insight**

107 One of our most important findings: **Training is working. It is only that Q3_difference**
108 **fails to capture it.**

- 109 • Smart systems as DV → $R^2 = .517$
- 110 • Q3_difference as DV → $R^2 = .136$

111 With the same IVs, simply changing the DV makes a **fourfold** difference in explanatory power. This
112 is the “dose-response separation” phenomenon that converged through 5-methodology triangulation
113 (OLS, PSM, DID, CLPM, RSM).

114 The AI was unable even to pose this question. It used only Q3_difference as the dependent variable
115 and never conceived the idea, “what if we change to a different DV?” This is the **absence of**
116 **domain knowledge**. Its ability to optimize within the data is outstanding, but the ability to
117 **imagine outside the data** is not yet there.

118 **Missing 4: Discussion**

119 The heart of a social-science paper is the Discussion. “What does this finding add to existing
120 theory?” “What does it mean for practitioners?” “What are the policy implications?”

121 The AI’s Conclusion stopped at a summary of results plus a list of limitations. The **prescriptive**
122 **conclusion** our interpretation report reached—the word “tailored,” the argument that we must
123 diagnose a firm’s current position and design interventions that supplement the missing dimensions—
124 was absent from the AI.

125

126 Chapter 4. So What Does This Mean?

127 The Real Value of the AI-Scientist

128 The R^2 .002 \rightarrow .336 that the AI found in 20 minutes and \$8 is one piece of the conclusion we
 129 reached in 3 months with 36 models. About one-fifth of the whole. But **the direction was**
 130 **exactly right**.

131 What this means: 1. **Value as an accelerator of exploration:** A workflow is possible in which
 132 the AI first scans for heterogeneity through clustering, and the researcher then layers on theory and
 133 triangulates. 2. **Value as a hypothesis-generation tool:** The AI's finding that "it works only
 134 in high readiness" provokes the next question, "why?" Answering that question with P-E Fit is the
 135 human's part. 3. **Value as a pattern-confirmation tool:** The AI independently confirms the
 136 conclusion we reached through 36 analyses. This is a kind of **methodological triangulation**.

137 The Limits of the AI-Scientist

- 138 1. **Absence of theory:** It finds patterns but cannot explain them.
- 139 2. **Absence of DV imagination:** It optimizes only within the given DV.
- 140 3. **Absence of literature:** There is no dialogue with prior research.
- 141 4. **Absence of Discussion:** It cannot answer "so what?"
- 142 5. **Absence of triangulation:** It relies on a single methodology.

143 All five of these are problems of **domain knowledge**. Its capacity for statistical pattern recognition
 144 already approaches or exceeds the human (20 minutes vs 3 months). What it lacks is knowing
 145 "what matters in this field."

146 The Next Step: AI-Scholar

147 Filling this gap is the development goal of the AI-Scholar.

- 148 • Theory Engine: theory DB + pattern-theory matching
- 149 • Literature Engine: SSCI literature search + citation generation
- 150 • Analysis Engine: automated triangulation (QCA + LPA + RSM + SEM)
- 151 • Writing Engine: a Theory \rightarrow Hypotheses \rightarrow Discussion structure

152 What the AI-Scientist demonstrated is **possibility**. In 20 minutes and \$8, it can make a finding
 153 that points in the right direction. Layer theory and context onto that, and the automatic generation
 154 of a social-science paper is not impossible.

156 **Closing: The Machine's Eye and the Researcher's Eye**

157 The AI saw $R^2 .002 \rightarrow .336$. We saw “when awareness and infrastructure meet.”

158 The AI wrote “dramatic heterogeneity.” We arrived at the word “tailored.”

159 The AI found the patterns. We named them, gave them meaning, and proposed what to do next.

160 Both of us looked at the same data, but **we look at different depths**. The AI's depth is
161 deepening fast. Yet the questions “why?” and “so what?” still belong to the human domain.

162 This preprint is the first social-science test of the AI-Scientist, and the starting point for developing
163 the AI-Scholar. How the machine's eye and the researcher's eye can collaborate is what we must
164 now design.