



Comparison of Freshman Baseline with First Year Seminar Assessment Results

Fall Semester 2024

Summer Assessment Team Members: Marie Archambault, Clinton Brown, Robert Ellison, Victor Fet, Marty Laubach, Leslie Dawn Quick, and Anita Walz

Summer Assessment Support Staff: Mary Beth Reynolds, Adam Russell, and Diana Adams

Executive Summary

Background

Recommendations from the 2024 Assessment Team

The Summer Assessment Team made the following recommendations:

1. That the FYS faculty consider some standardization regarding lessons on critical thinking, exposing students to the type of critical thinking and problem-solving on which the FYS final exam is based. For example, we recommend that, when presented with a problem or issue to address, students be given more experience examining evidence, considering multiple viewpoints, interrogating their own assumptions and biases in arriving at recommendations to address the issue or solutions to the problem. They should also consider possible consequences of their proposed recommendation or solution. **This recommendation was implemented in most FYS sections in 2024-2025.**
2. That FYS faculty work closely with the online Design Center and with University Libraries to ensure that source links to final exam documents are functional during that critical period of the semester. **One of our team members, who teaches large sections of FYS, reported no issues with final exam documents during academic year 2024-2025.**
3. That we again share the rubric we currently use with FYS faculty so that they are aware of the university's expectations for student performance on the exam. **To our knowledge, this was not done.**

Procedures for the 2024 Assessment

General Procedures

In August 2024, 1,545 incoming freshmen at Marshall University appeared to have uploaded baseline assessments into Blackboard as part of their assignments for Freshman First Class (UNI 100). These assessments required students to analyze and evaluate information, solve problems, and write effectively. These skills are aligned to three of Marshall University's outcomes; *Information Literacy*, *Inquiry-Based (Critical) Thinking*, and *Communication Fluency*. As part of Marshall's mandatory First Year Seminar in Critical Thinking (FYS), students completed assessments that mirrored those they finished as incoming freshmen, with 1,300 FYS assessments uploaded into Blackboard. To obtain a sample of matched pairs of baseline and FYS assessments, we began by comparing lists of all FYS and baseline artifacts uploaded to Blackboard during academic year 2024-2025 to determine which students submitted both baseline and FYS artifacts. We identified 463 potential matches and, from there, chose a random sample of 174 matched pairs. Each pair was further examined to ensure that the artifacts were uploaded. When this was not the case for either the baseline or FYS artifacts, that match was discarded, and another chosen until we had the desired 174 matched pairs. **Please note that, due to an issue with Blackboard, which we were not able to quickly fix, no FYS artifacts uploaded at the end of the spring 2025 semester were available for review. This greatly decreased the pairs from which we were able to pull our sample.**

In May 2025, a group of seven faculty representing three academic colleges (Liberal Arts, Science, and Business) evaluated the baseline/FYS sample using a rubric that allowed them to score each artifact across eight criteria (traits). These traits included information needed and source acknowledgment (*Information Literacy*), evidence, viewpoints, and recommendation/position (*Inquiry-Based [Critical] Thinking*), development, convention/format, and communication style (*Communication Fluency*). Given the ubiquitous availability of artificial intelligence (AI), they also decided to add a trait for suspected AI usage, bringing the total number of rubric traits to nine. This project was coordinated by the Office of Assessment and Quality Initiatives.

Each assessment had two independent raters. Please see the supporting documentation that follows this summary for a detailed explanation of scoring procedures.

Results and Analysis

Comparison of Freshman Baseline to Results at the End of FYS

The baseline and FYS means (and standard deviations) for the students in the sample with scorable baseline and FYS exams are reported below. We note that, despite the time spent checking the artifacts before scoring began, reviewers were either not able to access five uploaded artifacts or the artifacts that were accessed were missing sections B and C, which we use to assess these artifacts. Four additional artifacts were missing section C of the assignment, which covered all traits except *Information Needed*. This left us with 165 matched pairs scored for all traits, and 169 scored only for the first trait – *Information Needed*. We conducted *paired-samples t-tests* using adjusted alpha levels to control for Type I error (.025 for *Information literacy*), (.017 for *Inquiry-Based [Critical] Thinking*), and (.017 for *Communication Fluency*). Results showed significant differences between baseline and FYS results for all traits of each learning outcome. These results are shown in the table below. We further note that *Communication Fluency* is not an outcome of FYS.

Outcome	Trait	Baseline Mean (SD)	FYS Mean (SD)	Statistical Significance
Information Literacy	Information Needed	1.988 (0.6314)	2.387 (0.6349)	$t(168) = -6.591$, $p < .001$
	Source Acknowledgment	1.779 (0.6770)	2.676 (0.8093)	$t(164) = -11.548$, $p < .001$
Inquiry-Based (Critical) Thinking	Evidence	1.757 (0.6594)	2.467 (0.6855)	$t(164) = -10.581$, $p < .001$
	Viewpoints	1.642 (0.5348)	2.085 (0.4640)	$t(164) = -8.749$, $p < .001$
	Recommendation/Position	1.958 (0.6839)	2.503 (0.7017)	$t(164) = -8.246$, $p < .001$
Communication Fluency	Development	1.894 (0.7110)	2.642 (0.7567)	$t(164) = -10.469$, $p < .001$
	Convention/Format	1.961 (0.8002)	2.761 (0.7563)	$t(164) = -10.052$, $p < .001$
	Communication Style	2.415 (0.6064)	2.839 (0.4586)	$t(164) = -8.950$, $p < .001$

A frequency analysis also showed the following increases in students scoring between 2.5 and 4.0 on the rubric between baseline and FYS. Please see the supporting documentation following this summary for additional information.

Outcome	Trait	Percentage Gain in Students Scoring 2.5 to 4.0 from Baseline to FYS
Information Literacy	Information Needed	27%
	Source Acknowledgment	51%
Inquiry-Based (Critical) Thinking	Evidence	38%
	Viewpoints	16%
	Recommendation/Position	31%
Communication Fluency	Development	37%
	Convention/Format	42%
	Communication Style	1%

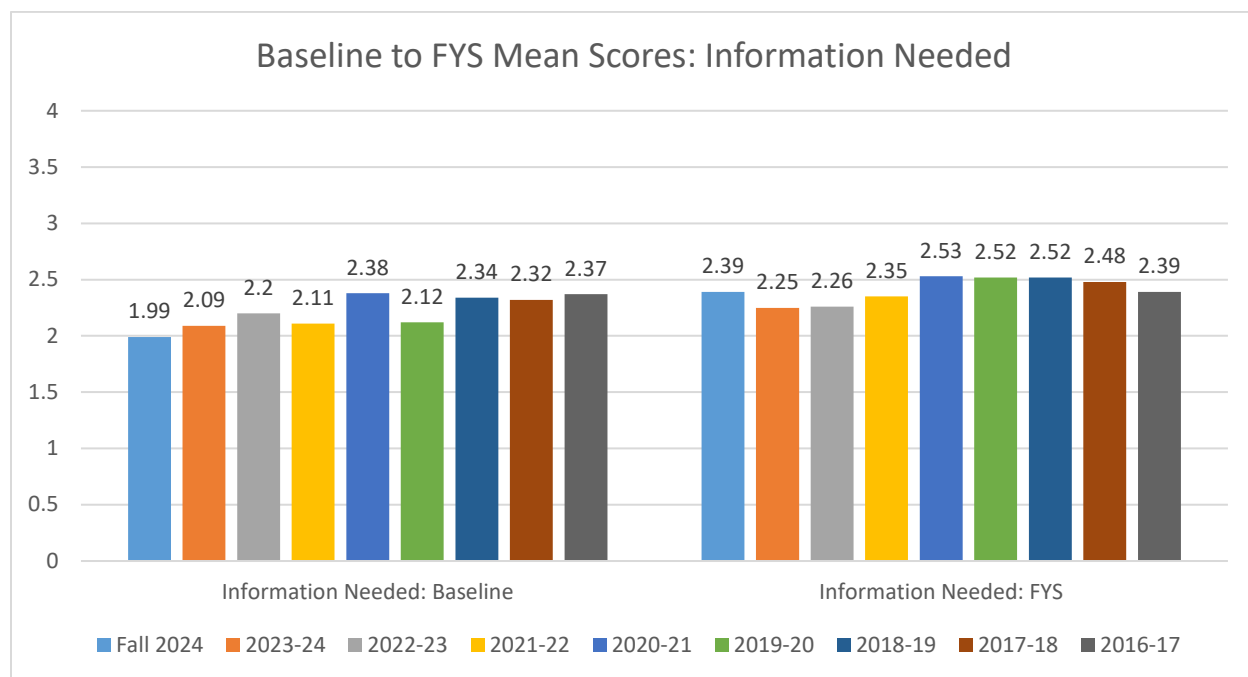
Since students enrolled in FYS in fall 2024 completed their responses to one of two possible scenarios, we further analyzed results based on scenario. Our sample included 100 students who completed the Flu Vaccine scenario and 74 who completed the GMO Food scenario. This year's results showed scores were higher on the Flu Vaccine scenario than for the GMO Food scenario across all traits of each outcome, with the differences reaching statistical significance for six of the eight traits (Source Acknowledgment for Information Literacy); (Evidence and Recommendations/Position for Inquiry-Based [critical] thinking), and for all traits aligned with communication fluency.

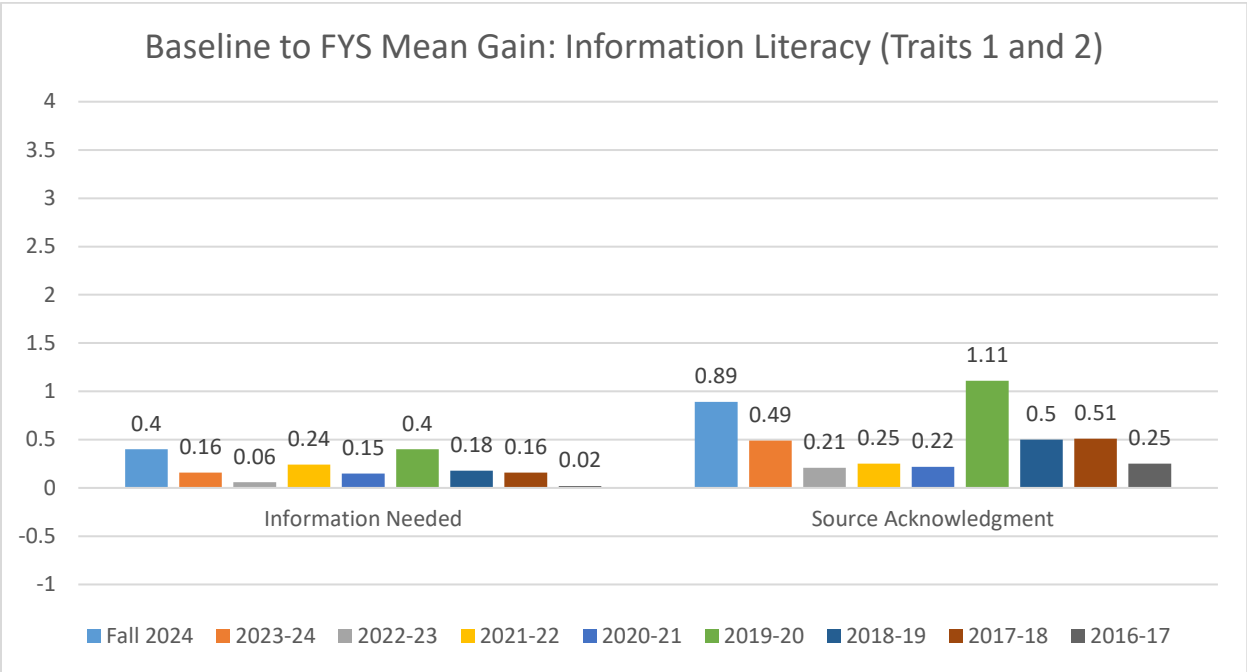
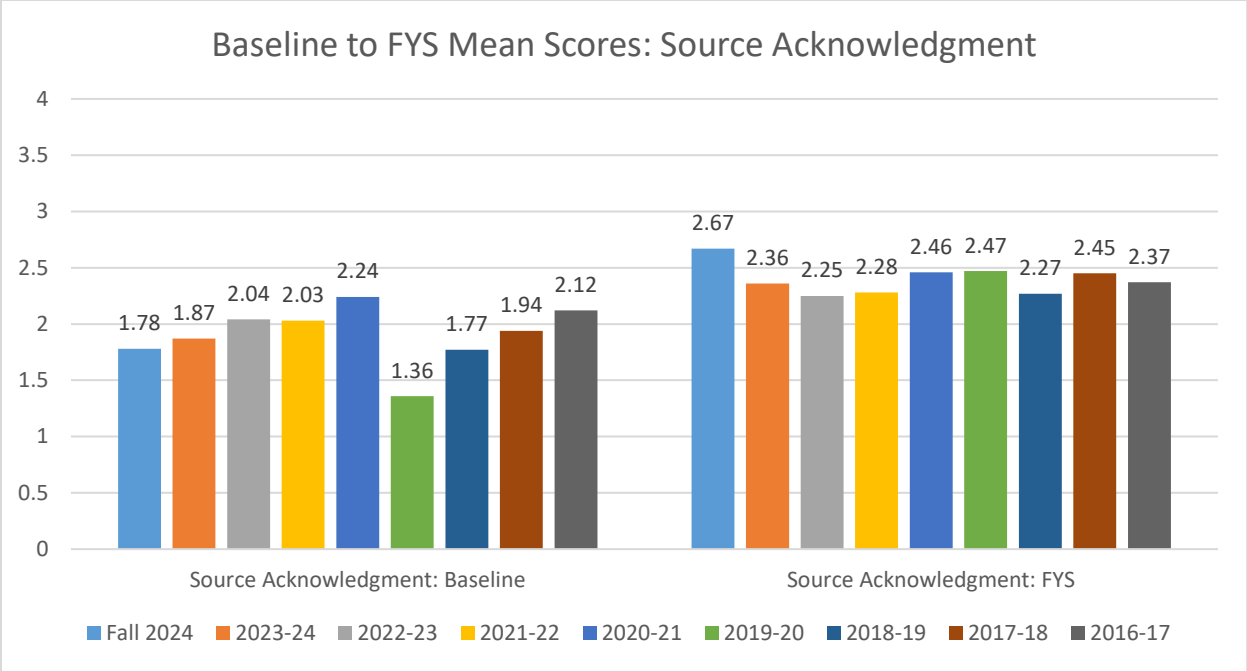
Conclusions

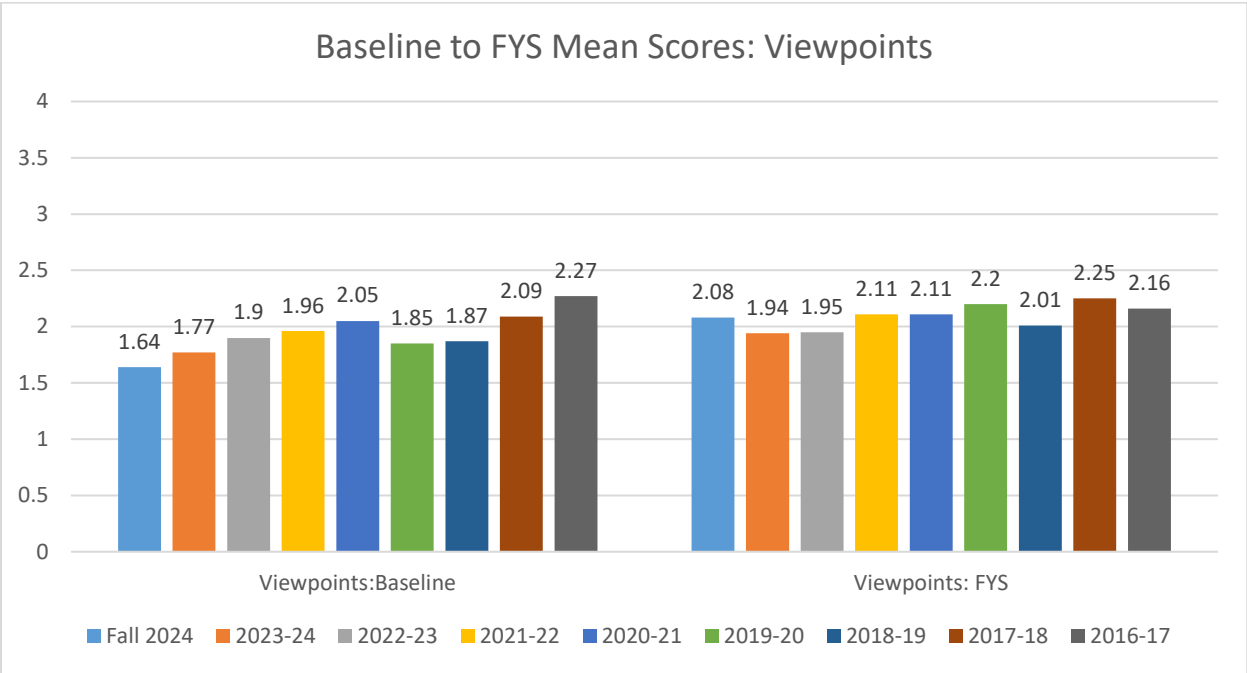
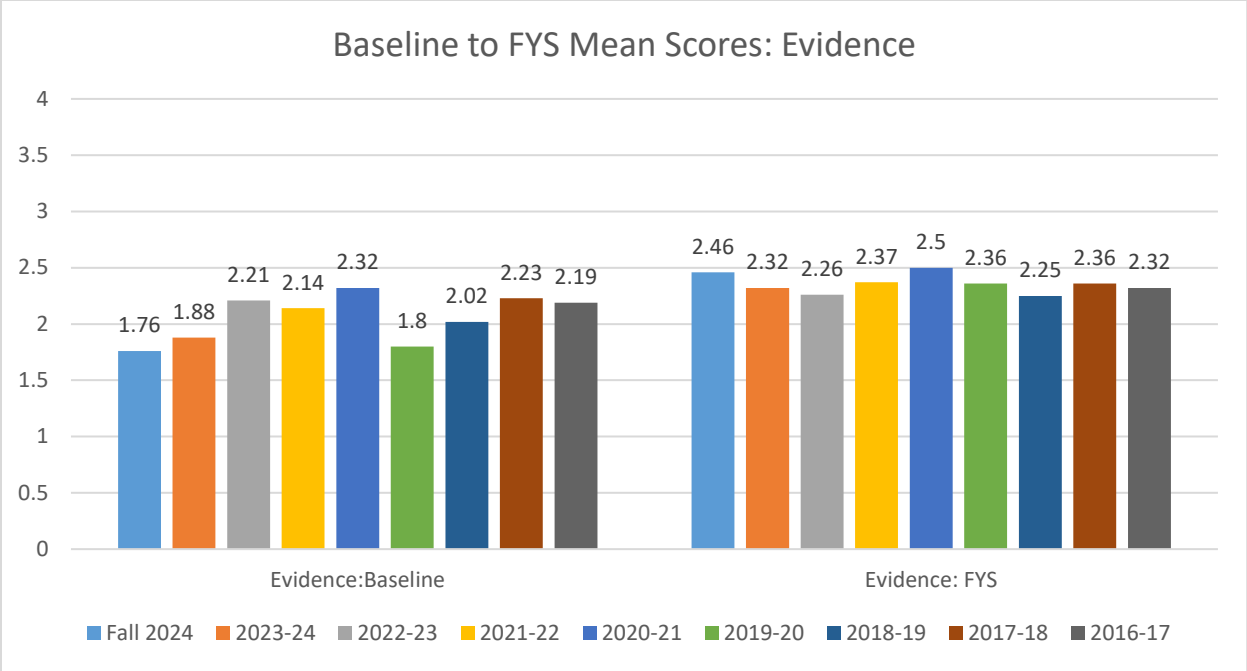
Although we have not performed statistical analyses to compare the results across years, we were concerned about our results in 2022-2023 because that was the only year we had not seen statistically significant improvement between baseline and FYS in at least some traits of *Critical Thinking* since we

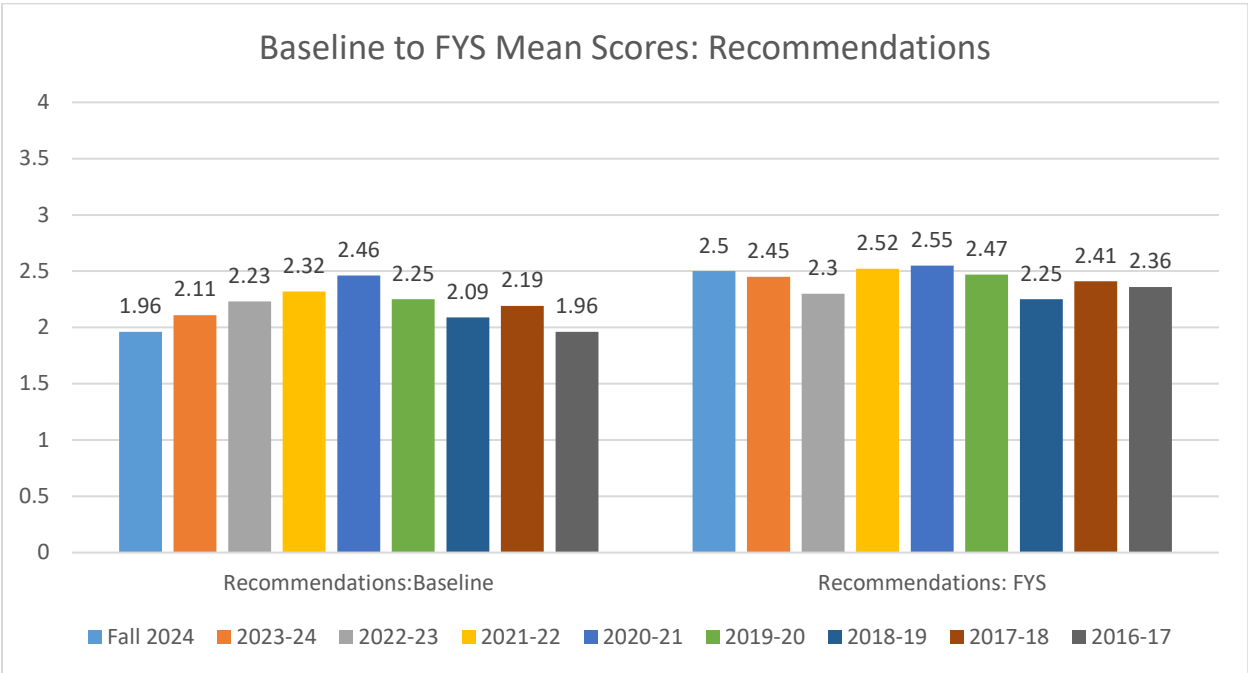
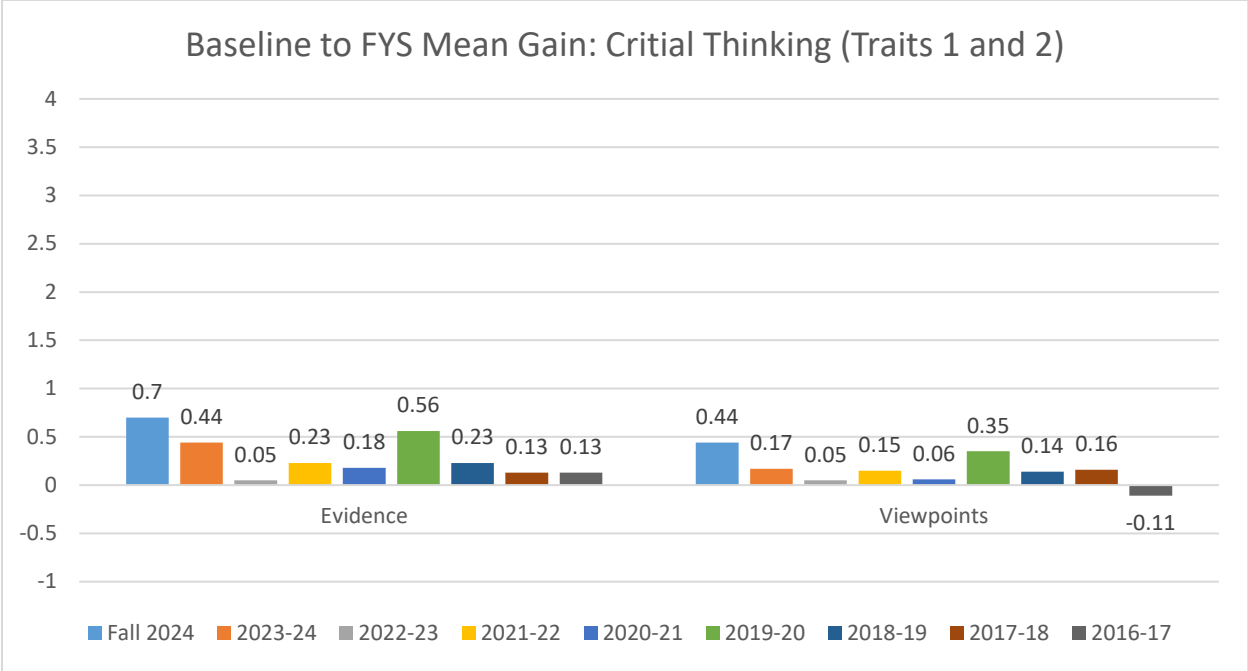
began analyzing student performance in 2013. After comparing trends for baseline and FYS means from 2016-2017 through 2022-2023, we concluded that the 2022-2023 results were not due to higher than usual baseline scores, but rather to lower FYS scores than in past years.

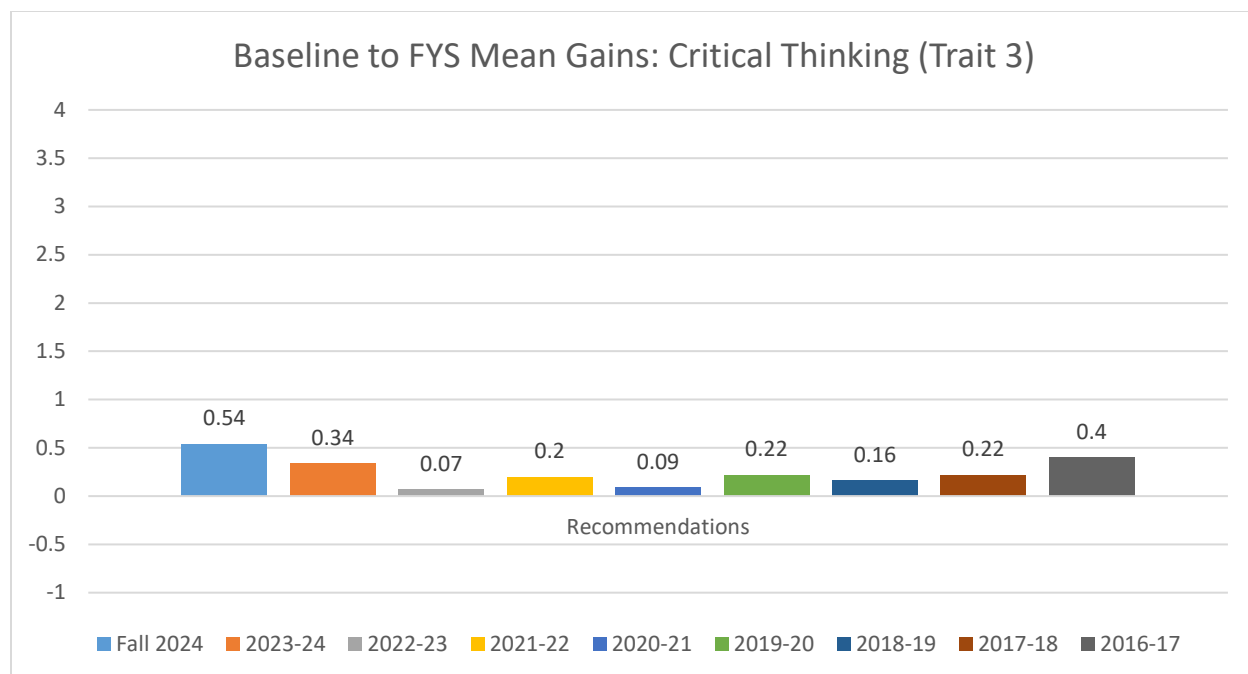
As noted, 2023-2024 and fall 2024's results did show significant differences for all traits of the three outcomes (*Information Literacy*, *Critical Thinking*, and *Communication Fluency*) assessed. In reviewing these results, we noted that 2023-2024 and fall 2024 students scored lower on baseline and higher in most cases on FYS than students from 2022-2023. This led us to examine two metrics – 1) gain score for *Information Literacy* and *Critical Thinking* for our samples from 2016 to the present, and 2) Baseline and FYS means for the same period. This information is shown below.











Examination of the charts above show that, for the most part, mean scores at the end of FYS reach between 2.0 and 2.4 on a 4-point rubric scale, with students making gains between baseline and FYS for all rubric traits except *Critical Thinking (viewpoints)* in one out of the nine years examined.

Recommendations from the 2025 Assessment Team

The Summer Assessment Team made the following recommendations:

1. Given the consistent differences between student performance on the two scenarios from fall 2024, the Summer Assessment Team (SAT) suggests that the baseline/FYS team standardize the deliverables requested of the students, e.g., a letter or memorandum requires the students to address the elements of the rubric at a higher level than does an op-ed.
2. Before launching a revised general education curriculum, the SAT suggests that the General Education Task Force meet with the University Assessment Committee to determine a workable plan for general education assessment. This plan should include a method for faculty to improve curriculum and pedagogy based on the results of general education assessment.



Supporting Documentation



Comparison of Freshman Baseline and First-Year Seminar (FYS) Assessments

Fall Semester 2024

Review Procedures

- One hundred seventy-four (174) FYS critical thinking artifacts were matched with 174 baseline critical thinking artifacts. This number represented 13% of the 1,300 FYS artifacts and 11% of the 1,529 baseline artifacts uploaded to Blackboard.
- During the evaluation we discovered that four of the baseline artifacts included only the information linked to the *Information Needed* part of the rubric and an additional five included no information aligned to the rubric or were not able to be accessed. This left reviewers with 169 matched pairs for the *Information Needed* rubric trait and 165 matched pairs for all other traits.

Review Procedures Continued

- Each assessment had two independent raters and scores were determined in the following manner:
 - If raters assigned the same score, that became the score for the artifact.
 - If raters' scores differed by one point, e.g., Rater 1 assigned a score of 1 and Rater 2 a score of 2, the final score was the mean, i.e., 1.5.
 - If raters' scores differed by more than one point, e.g., Rater 1 assigned a score of 1 and Rater 2 a score of 3, the raters met to discuss the rationale for their scores to see if they could agree on a score or, at minimum, scores that differed by no more than one point.
 - If raters' scores differed by more than one point and, after discussion, they were not able to resolve the differences, a third rater was assigned to review the assessment. (For this review, all raters were able to reconcile disagreements, so third raters were not needed).

Interrater Reliability

- We conducted interrater reliability analyses using the Cohen's Kappa statistical procedure. In so doing, we used the following rules, similar to those suggested by Stellmack, Kohneim-Kalkstein, Manor, Massey, & Schmitz (2009):
 - Since our scoring procedure was to average final scores between two raters when scores differed by only one point, we used that averaged score (e.g., 1.5) as the score for both raters, counting it as an agreement in the interrater reliability analysis.
 - For scores that were two or more points apart, the original score of each reviewer was used in the analysis. Therefore, these scores were counted as disagreements.

Rubric Used for Scoring

Baseline/FYS Assessment Rubric – Summer 2023 – updated 5-8-2023

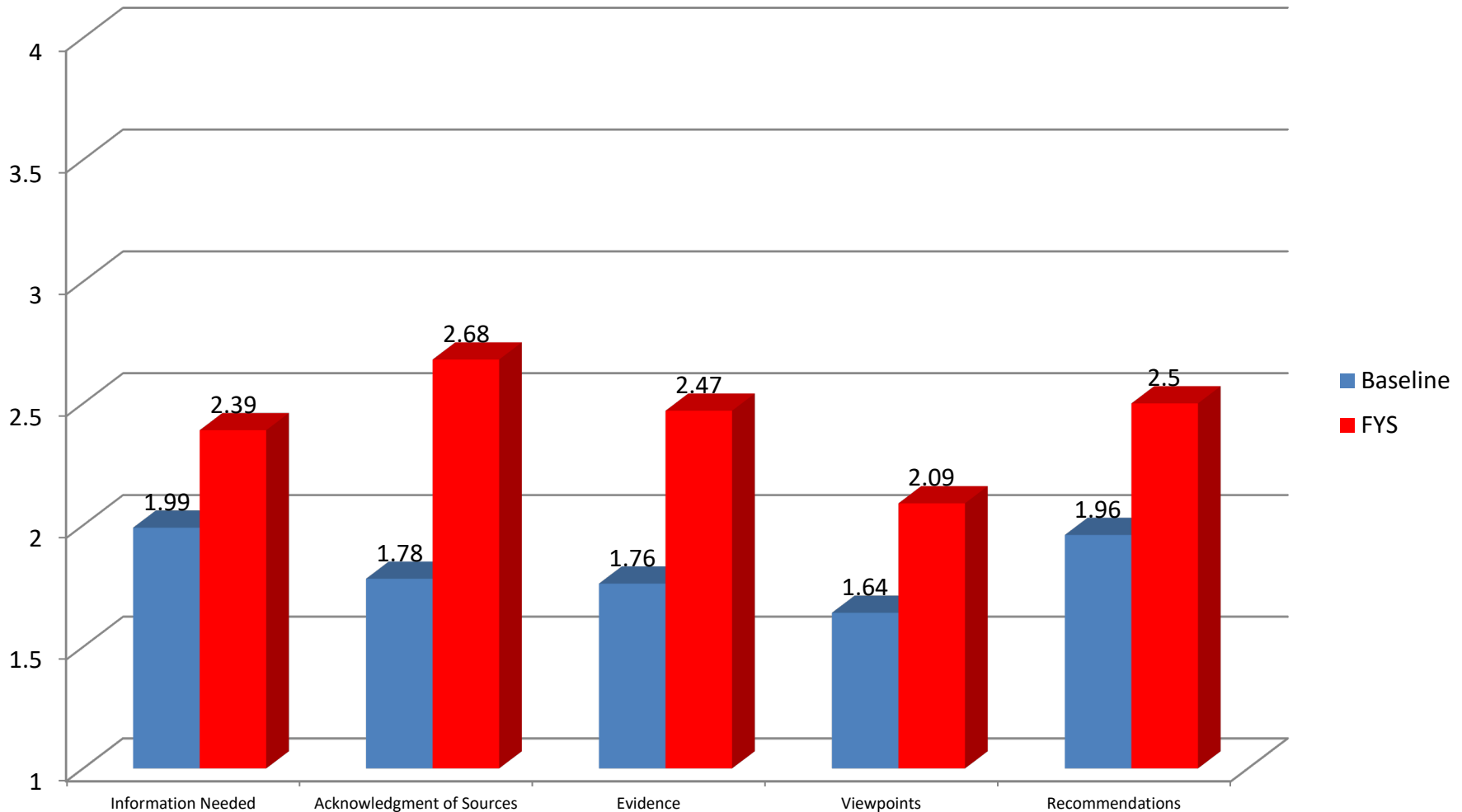
Outcomes	Traits	Performance Levels			
		1	2	3	4
Information Literacy	Information Needed	Does not acknowledge or assess the need for more information.	Acknowledges the need for more information but does not identify research methods/sources (or those identified are not feasible) that would address unanswered questions.	Assesses the need for more information and recommends general research methods/sources (that are feasible) that would address some unanswered questions.	Assesses the need for more information and recommends specific research methods/sources (that are feasible) that would address most unanswered questions.
	Source Acknowledgment	Fails to acknowledge sources from the DL.	Indirectly/vaguely acknowledges sources of information from the DL.	Clearly acknowledges relevant sources of information from the DL.	Integrates relevant information from the DL. Acknowledges sources used.
Inquiry-Based Thinking	Evidence	Disregards or misunderstands evidence from the DL.	Insufficient evidence is taken from sources (e.g., only one or two pieces of evidence) in the DL or evidence is used without appropriate interpretation/evaluation (i.e., poor job).	Evidence is taken from relevant and valid sources in the DL with some interpretation/evaluation, but not enough to develop a coherent analysis or synthesis (i.e., adequate job).	Evidence is taken from relevant and valid sources in the DL with enough interpretation/evaluation to develop a coherent analysis or synthesis (i.e., good/excellent job).
	Viewpoints	Ignores viewpoints expressed in the DL.	Viewpoints expressed in the DL are taken as mostly fact, with little or no question.	Questions some viewpoints expressed in the DL.	Thoroughly questions and evaluates viewpoints expressed in the DL.
	Recommendation/Position	<u>Either</u> does not make a recommendation (take a position) <u>or</u> makes a recommendation (takes a position), but does not justify it in any way.	Recommendation/position is justified, but does not acknowledge different sides of the issue.	Recommendation/position is justified and takes into account different sides/complexities of the issue.	Recommendation/position takes into account the complexities of the issue. Any limits to the recommendation are acknowledged.
Communication Fluency	Development	Shows little or no evidence of developing their ideas.	Shows some development of ideas.	Shows a strong, but perhaps somewhat incomplete, development of ideas.	Produces a document in which the ideas have been fully developed.
	Convention/Format	Demonstrates minimal attention to basic organization and presentation and stylistic conventions.	Demonstrates some awareness of basic organization, content, and presentation and stylistic conventions.	Demonstrates consistent use of important conventions particular to a specific writing task, including organization, content, presentation, and stylistic choices.	Demonstrates detailed attention to and successful execution of a wide range of conventions particular to a specific writing task including organization, content, presentation, formatting, and stylistic choices.
	Communication Style	Uses language that impedes meaning because of errors in usage/mechanics.	Uses language that generally conveys meaning to readers, although errors in usage/mechanics impedes smooth reading of the document.	Uses straightforward language that conveys meaning to readers. The language in the document has few errors.	Uses language that skillfully communicates meaning to readers with clarity and fluency, and is virtually error-free.

Freshman Baseline/FYS Comparisons

Mean Scores on a scale of 1 – 4, with 4 being the highest possible score

$n = 169$ matched pairs (Information Needed) and 165 matched pairs for all other traits.

Mean differences between baseline and FYS were statistically significant for all traits.



Freshman Baseline/FYS Comparisons

n = Baseline/FYS comparisons = 169 (Information Needed); 165 (all other traits)

Trait/ Performance Level	Info Needed (Baseline)	Info Needed (FYS)	Acknowledgment of Sources (Baseline)	Acknowledgment of Sources (FYS)
1.0	25 (15%)	10 (6%)	42 (26%)	16 (10%)
1.5 – 2.0	96 (57%)	67 (40%)	86 (52%)	27 (16%)
2.5 – 3.0	41 (24%)	79 (47%)	34 (21%)	83 (50%)
3.5 – 4.0	7 (4%)	13 (8%)	3 (2%)	39 (24%)
Totals	169	169	165	165

Freshman Baseline/FYS Comparisons

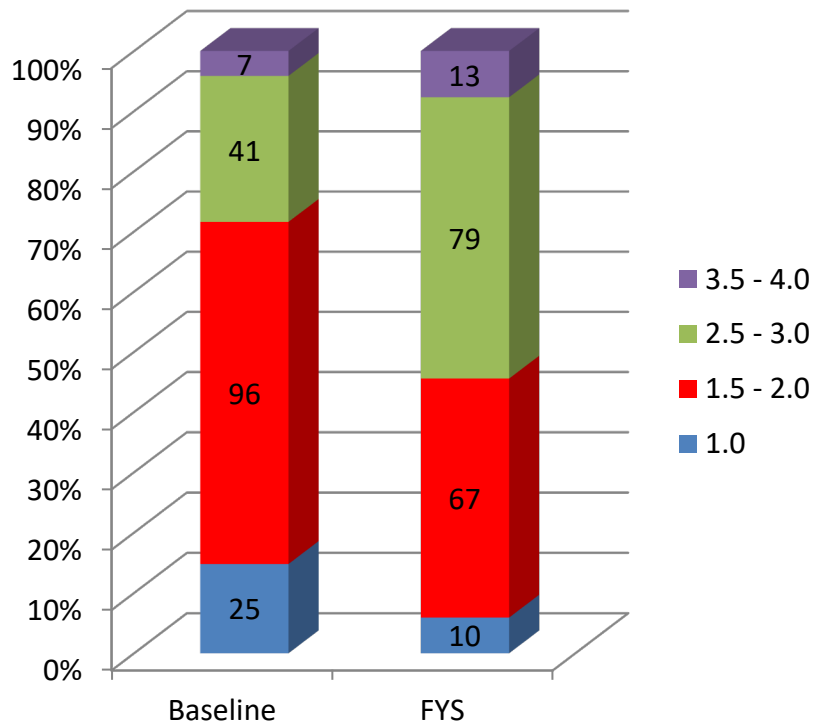
$n = 165$

Trait/ Performance Level	Evidence Baseline	Evidence FYS	Viewpoints Baseline	Viewpoints FYS	Recommendations Baseline	Recommendations FYS
1.0	47 (29%)	11 (7%)	51 (31%)	7 (4%)	37 (23%)	14 (8%)
1.5 – 2.0	76 (46%)	48 (29%)	98 (59%)	114 (69%)	68 (41%)	41 (25%)
2.5 – 3.0	41 (25%)	86 (52%)	16 (10%)	42 (25%)	58 (35%)	94 (57%)
3.5 – 4.0	1 (1%)	20 (12%)	0	2 (1%)	2 (1%)	16 (10%)
Totals	165	165	165	165	165	165

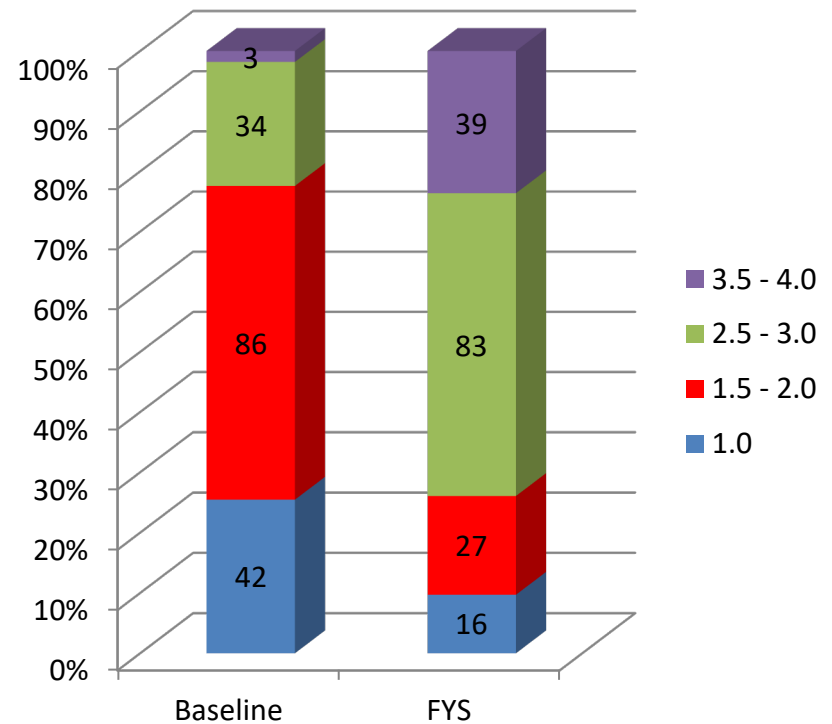
Freshman Baseline/FYS Comparisons

n = 169 matched pairs (Information Needed); 165 matched pairs-all other traits

Information Needed



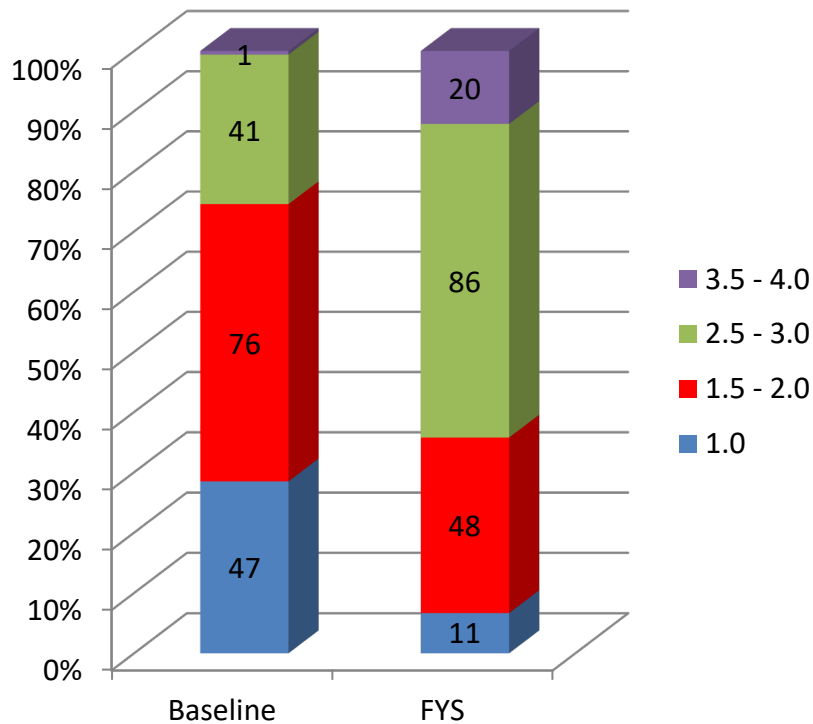
Acknowledgment of Sources



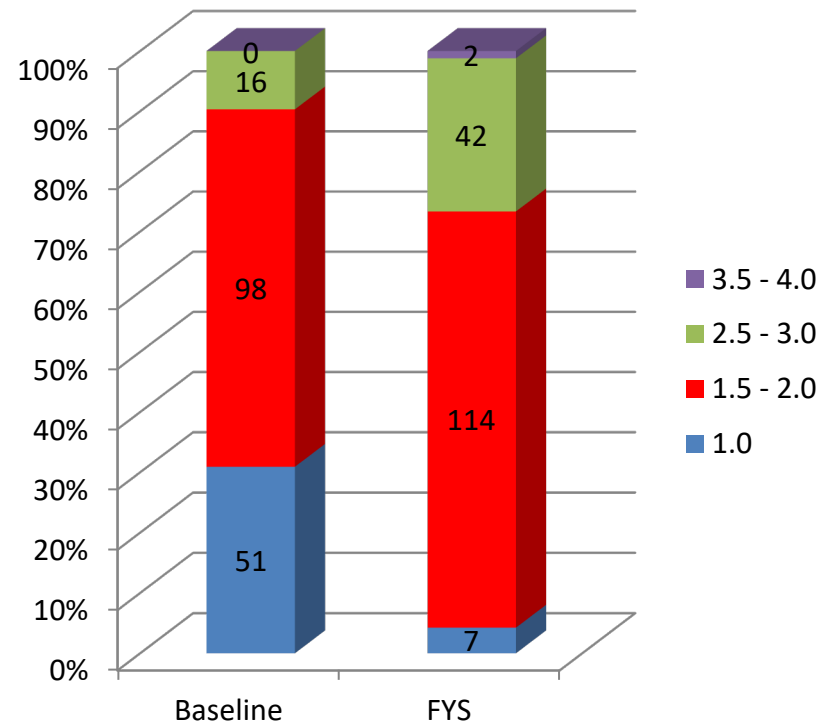
Freshman Baseline/FYS Comparisons

$n = 165$ matched pairs

Evidence



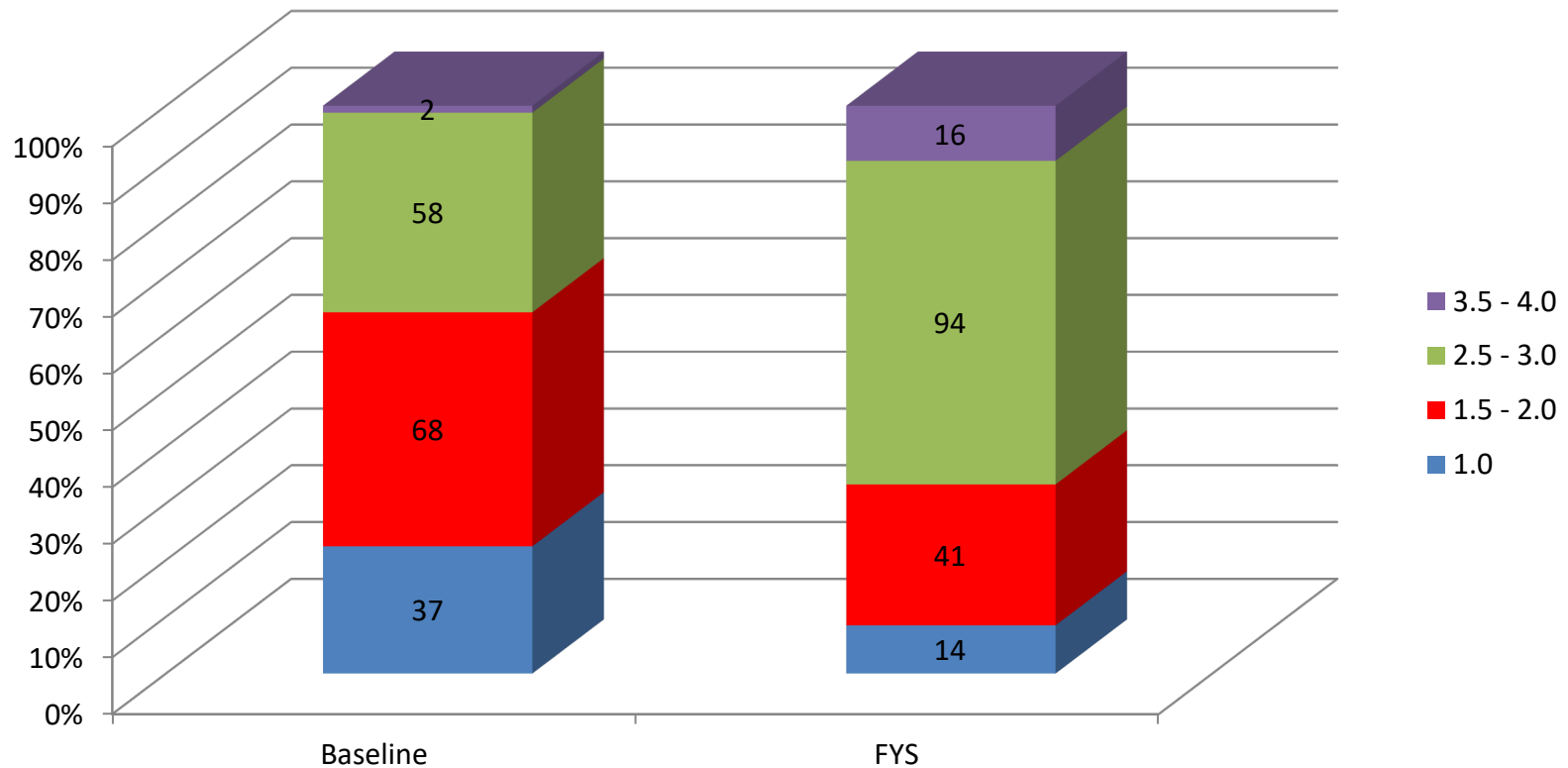
Viewpoints



Freshman Baseline/FYS Comparisons

$n = 165$ matched pairs

Recommendations



Baseline Inter-Rater Agreement Results

Includes 169 (Information Needed); 165 (all other traits)

Trait/ Agreement	Info Needed : Cohen's Kappa (Liberal) = .976	Acknowledgment of Sources: Cohen's Kappa (Liberal) = .992	Evidence: Cohen's Kappa (Liberal) = .961	Viewpoints: Cohen's Kappa (Liberal) = .975	Recommendations: Cohen's Kappa (Liberal) = .931
Agree on score	104 (62%)	95 (58%)	91 (55%)	116 (70%)	98 (59%)
Difference = 1 point	62 (37%)	69 (42%)	69 (42%)	46 (28%)	58 (35%)
Difference = 2 points	3 (2%)	1 (1%)	5 (3%)	3 (2%)	9 (5%)
Difference = 3 points	0	0	0	0	0
Total	169	165	165	165	165

FYS Inter-Rater Agreement Results

Includes all 174 FYS assessments scored

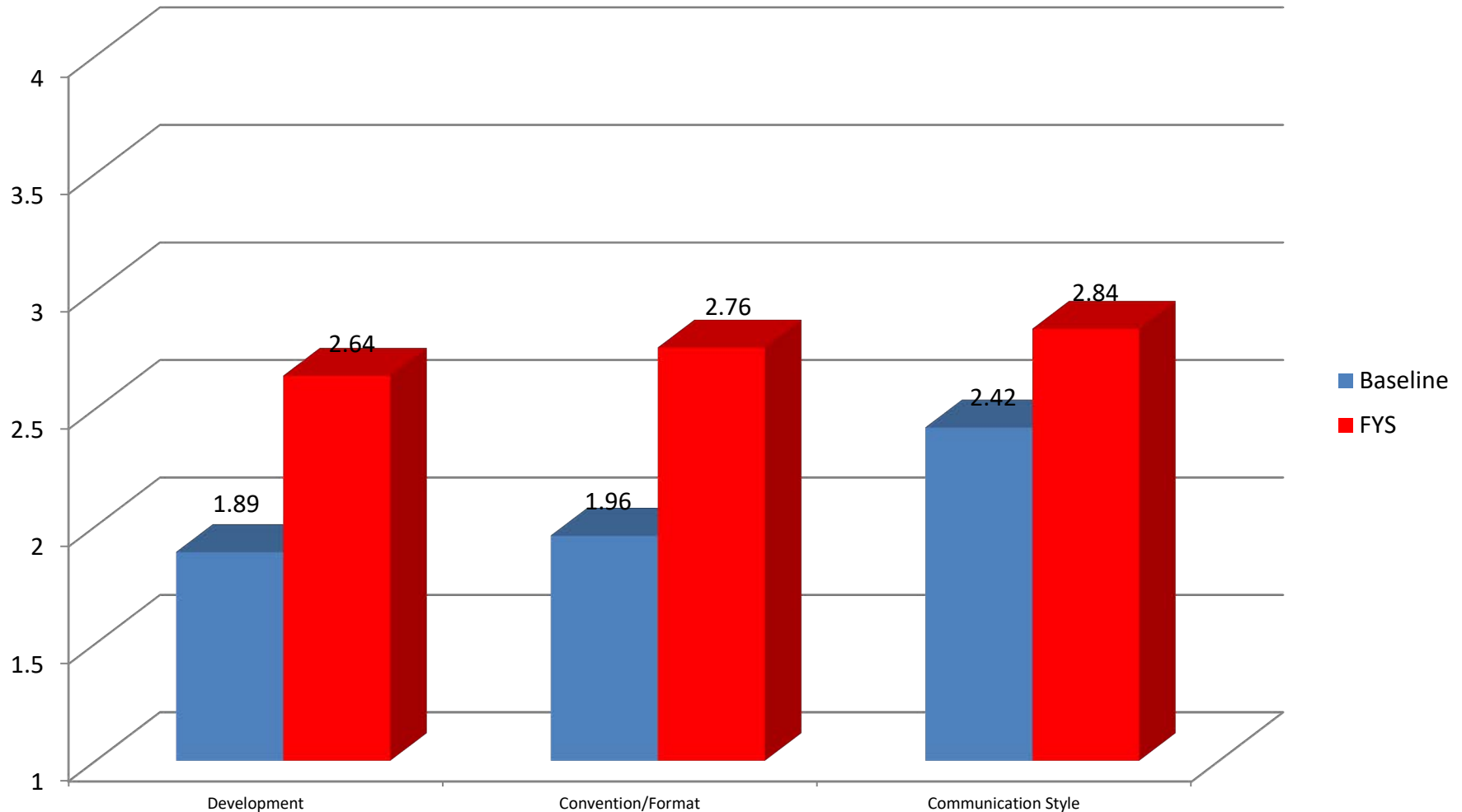
Trait/ Agreement	Info Needed : Cohen's Kappa (Liberal) = .947	Acknowledgment of Sources: Cohen's Kappa (Liberal) = .945	Evidence: Cohen's Kappa (Liberal) = .921	Viewpoints: Cohen's Kappa (Liberal) = .971	Recommendations: Cohen's Kappa (Liberal) = .901
Agree on score	109 (63%)	99 (57%)	89 (51%)	122 (70%)	108 (62%)
Difference = 1 point	58 (33%)	68 (39%)	74 (43%)	49 (28%)	53 (30%)
Difference = 2 points	7 (4%)	7 (4%)	11 (6%)	3 (2%)	13 (7%)
Difference = 3 points	0	0	0	0	0
Total	174	174	174	174	174

Freshman Baseline/FYS Comparisons

Mean Scores on a scale of 1 – 4, with 4 being the highest possible score

$n = 172$ matched pairs

Mean differences between baseline and FYS were statistically significant for *all traits*.



Freshman Baseline/FYS Comparisons

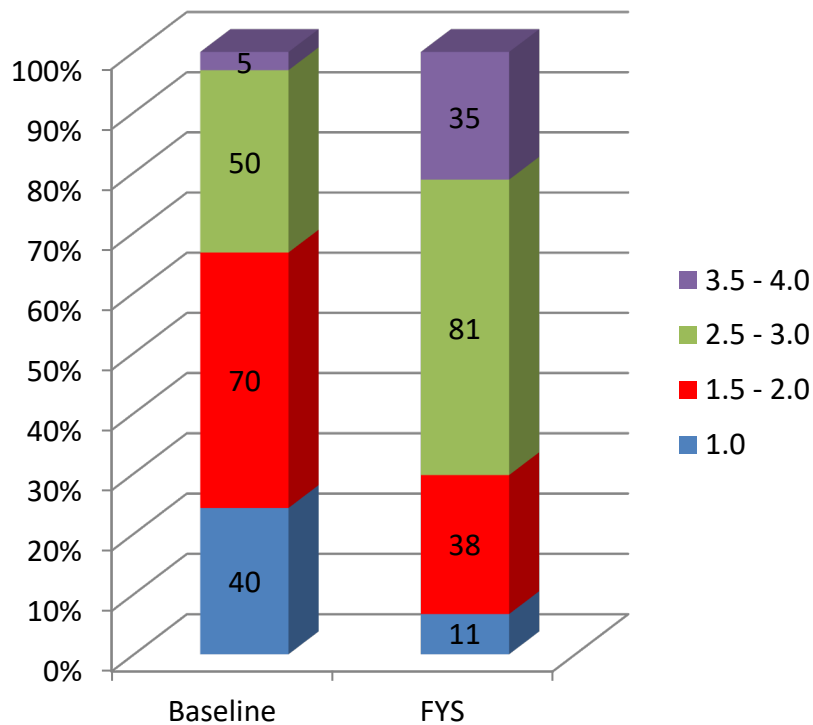
$n = 165$

Trait/ Performance Level	Development Baseline	Development FYS	Convention/ Format Baseline	Convention/ Format FYS	Communication Style Baseline	Communication Style FYS
1.0	40 (24%)	11 (7%)	40 (24%)	11 (7%)	7 (4%)	2 (1%)
1.5 – 2.0	70 (42%)	38 (23%)	64 (39%)	24 (15%)	51 (31%)	54 (33%)
2.5 – 3.0	50 (30%)	81 (49%)	48 (29%)	81 (49%)	99 (60%)	109 (66%)
3.5 – 4.0	5 (3%)	35 (21%)	13 (8%)	49 (30%)	8 (5%)	0
Totals	165	165	165	165	165	165

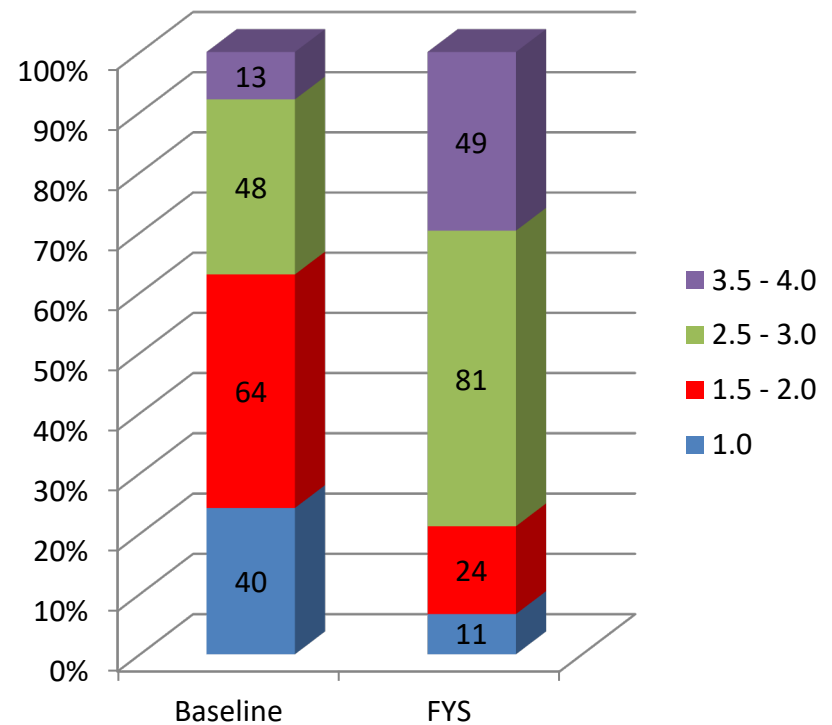
Freshman Baseline/FYS Comparisons

$n = 165$ matched pairs

Development



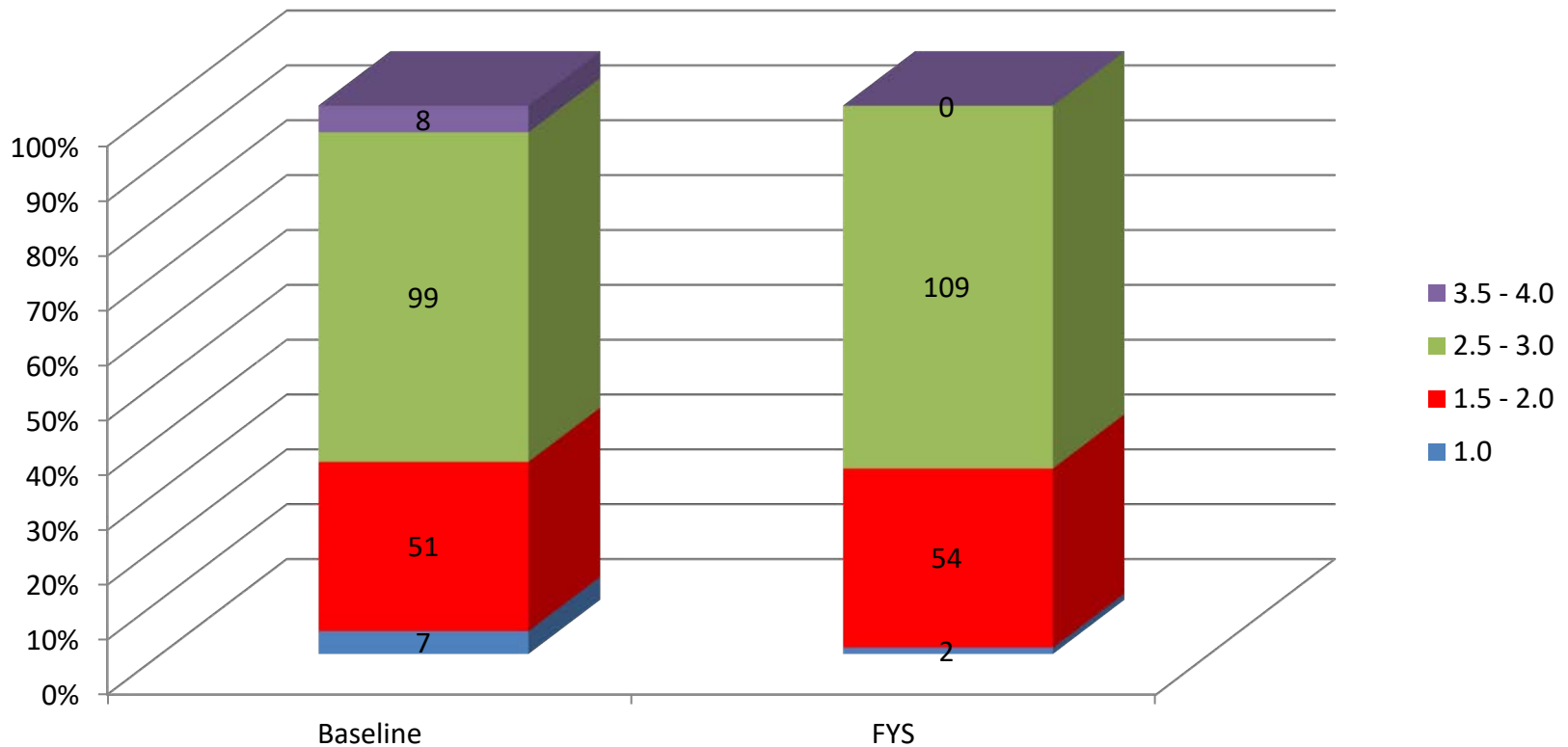
Convention/Format



Freshman Baseline/FYS Comparisons

$n = 165$ matched pairs

Communication Style



Baseline Inter-Rater Agreement Results

Includes 165 baseline assessments scored

Trait/ Agreement	Development: Cohen's Kappa (Liberal) = .962	Convention/Format: Cohen's Kappa (Liberal) = .941	Communication Style: Cohen's Kappa (Liberal) = .937
Agree on score	89 (54%)	83 (50%)	82 (50%)
Difference = 1 point	71 (43%)	74 (45%)	75 (45%)
Difference = 2 points	5 (3%)	8 (5%)	8 (5%)
Difference = 3 points	0	0	0
Total	165	165	165

FYS Inter-Rater Agreement Results

Includes all 174 baseline assessments scored

Trait/ Agreement	Development: Cohen's Kappa (Liberal) = .887	Convention/Format: Cohen's Kappa (Liberal) = .908	Communication Style: Cohen's Kappa (Liberal) = .931
Agree on score	91 (52%)	81 (47%)	101 (58%)
Difference = 1 point	67 (39%)	80 (46%)	65 (17%)
Difference = 2 points	16 (9%)	13 (7%)	8 (5%)
Difference = 3 points	0	0	0
Total	174	174	174



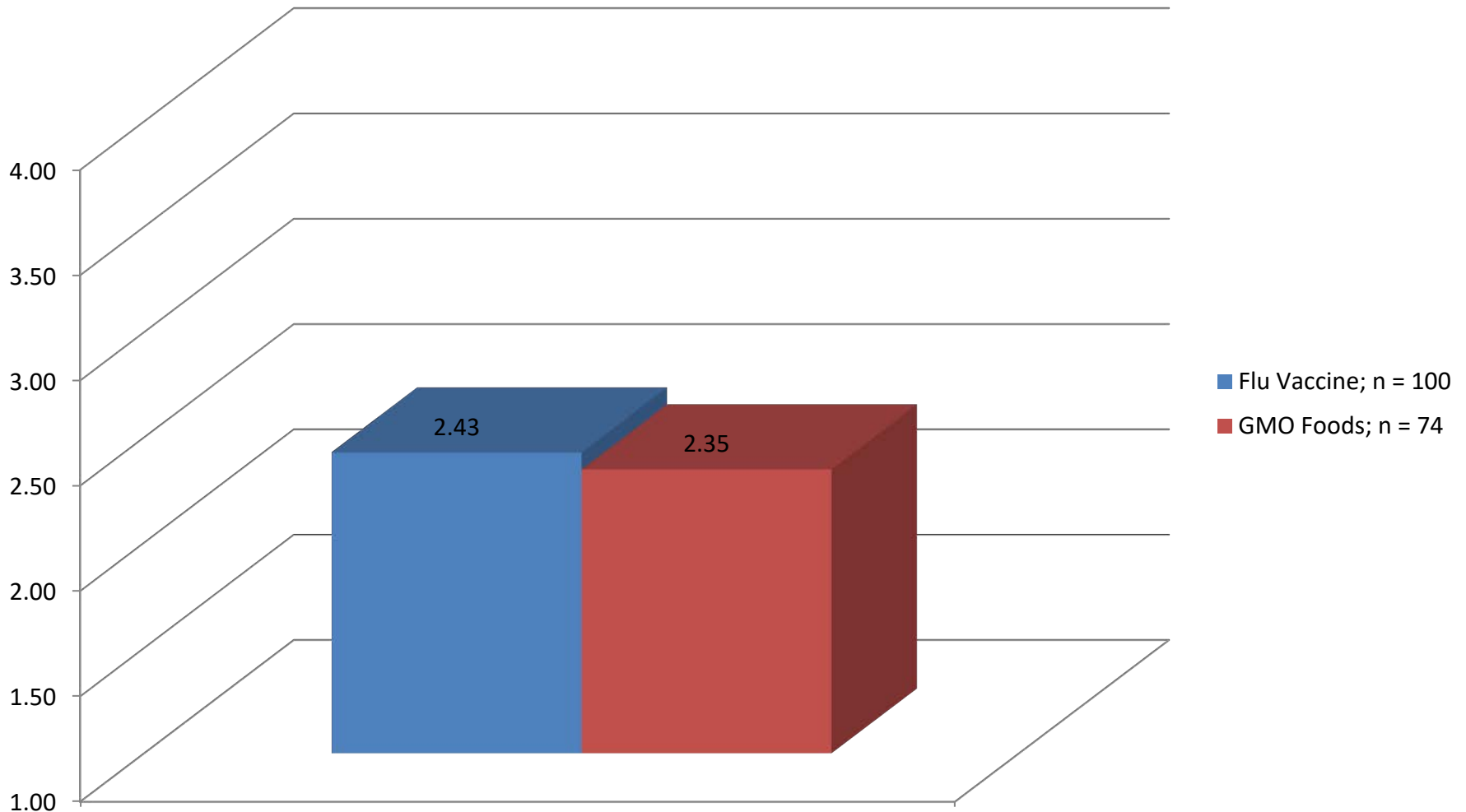
Comparison of FYS Results for Each Trait by Scenario

Fall 2024

FYS Comparisons by Scenario for IL: Information Needed

Mean Scores on a scale of 1 – 4, with 4 being the highest possible score

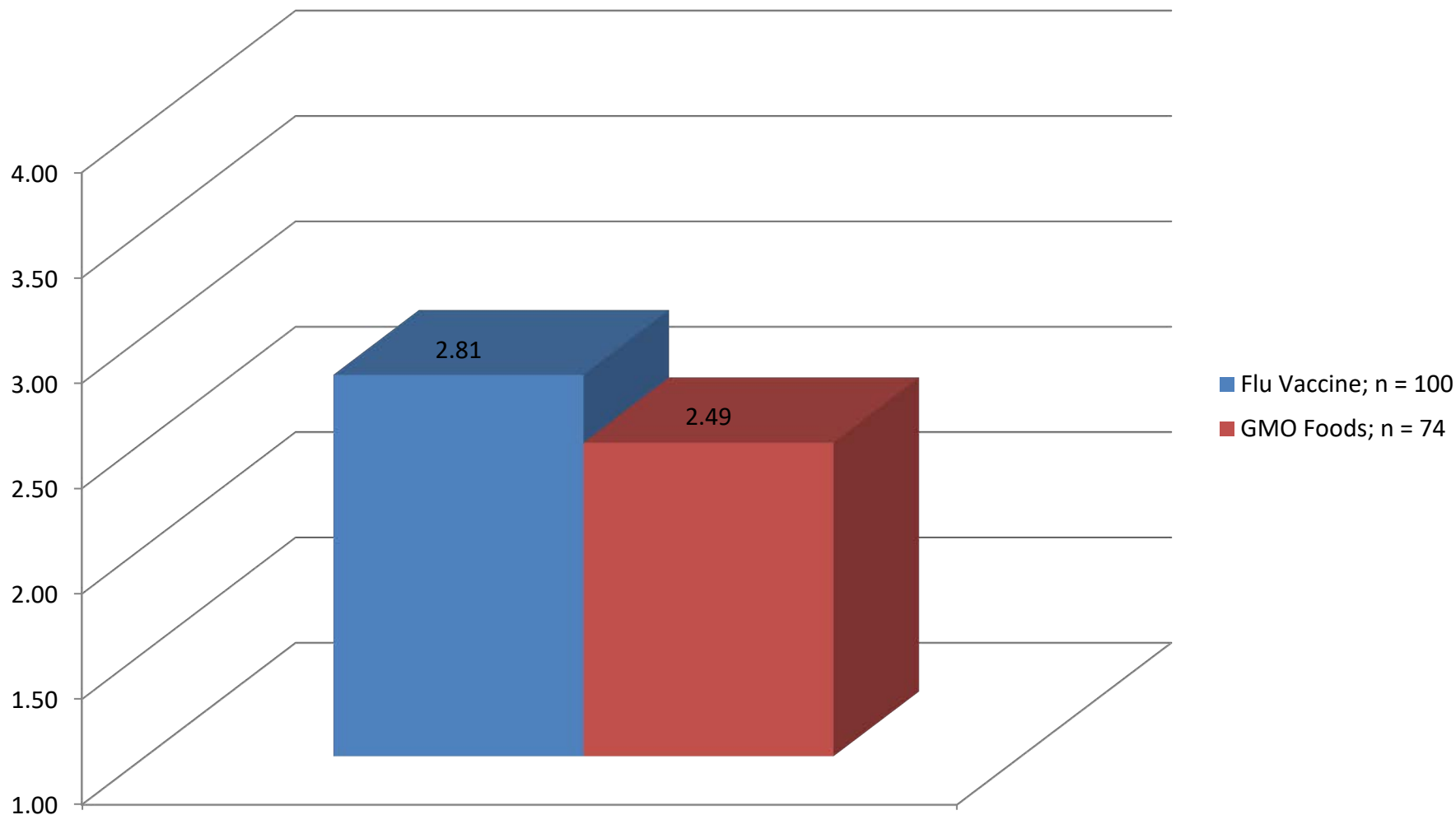
An independent samples t-test revealed no statistically significant differences between the means of these scenarios.



FYS Comparisons by Scenario for IL: : Source Acknowledgment

Mean Scores on a scale of 1 – 4, with 4 being the highest possible score

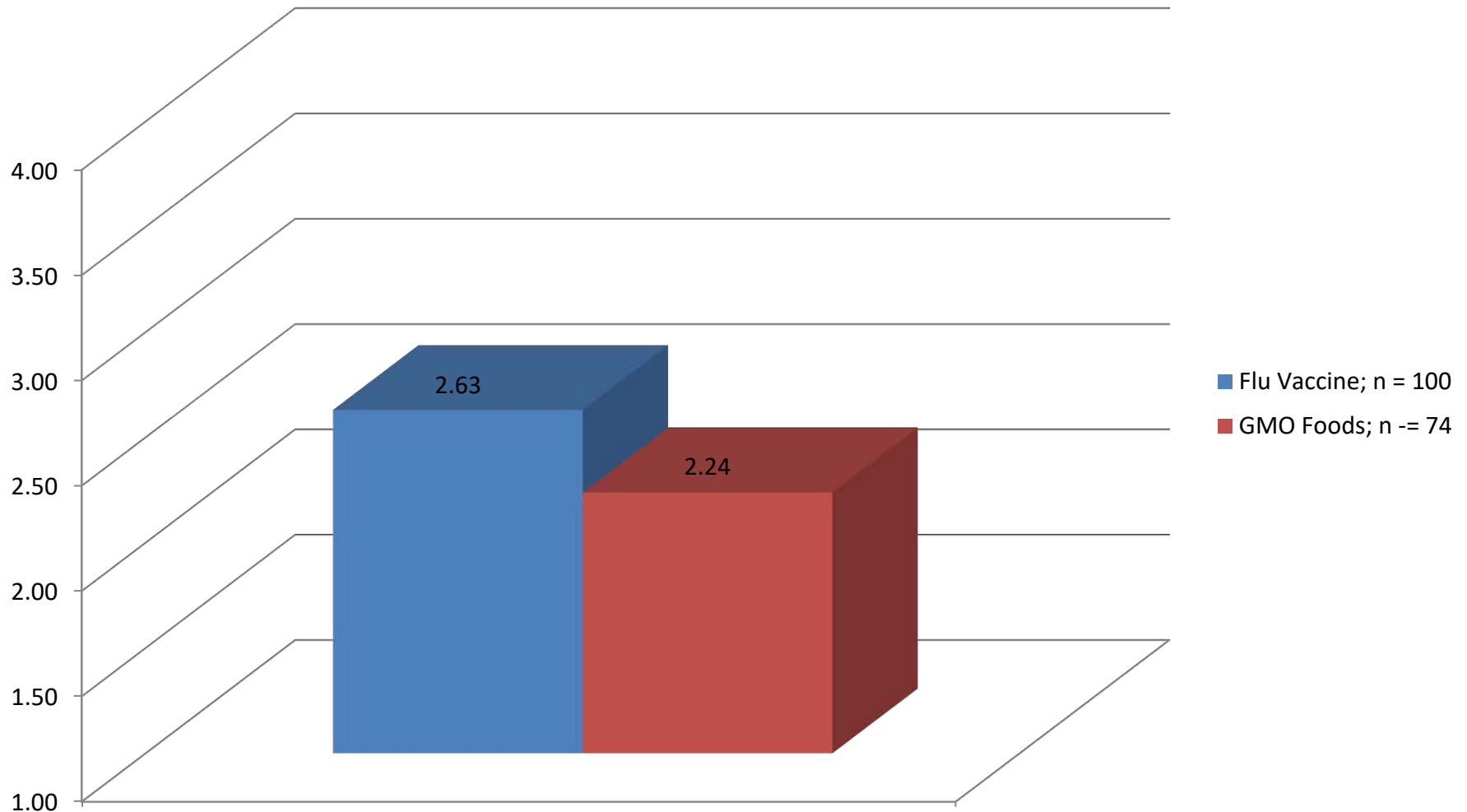
Using an adjusted alpha level of .025, an Independent Samples t-test revealed a statistically significant difference between the Flu Vaccine and GMO Foods Scenarios, $t(172) = 2.599, p = .010$.



FYS Comparisons by Scenario for CT: Evidence

Mean Scores on a scale of 1 – 4, with 4 being the highest possible score

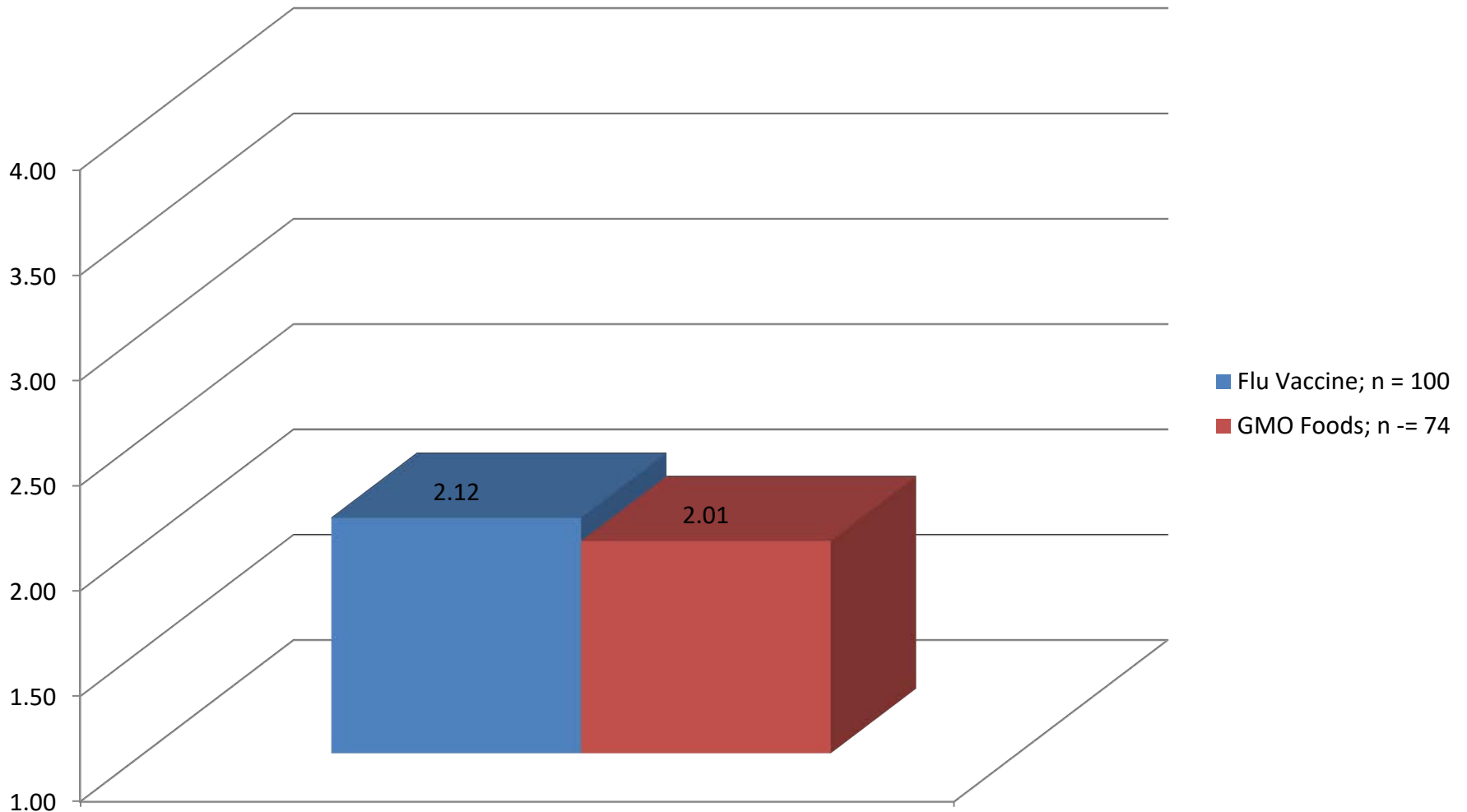
Using an adjusted alpha level of .017, an Independent Samples t-test revealed a statistically significant difference between the Flu Vaccine and GMO Foods Scenarios, $t(172) = 3.928, p < .001$.



FYS Comparisons by Scenario for CT: Viewpoints

Mean Scores on a scale of 1 – 4, with 4 being the highest possible score

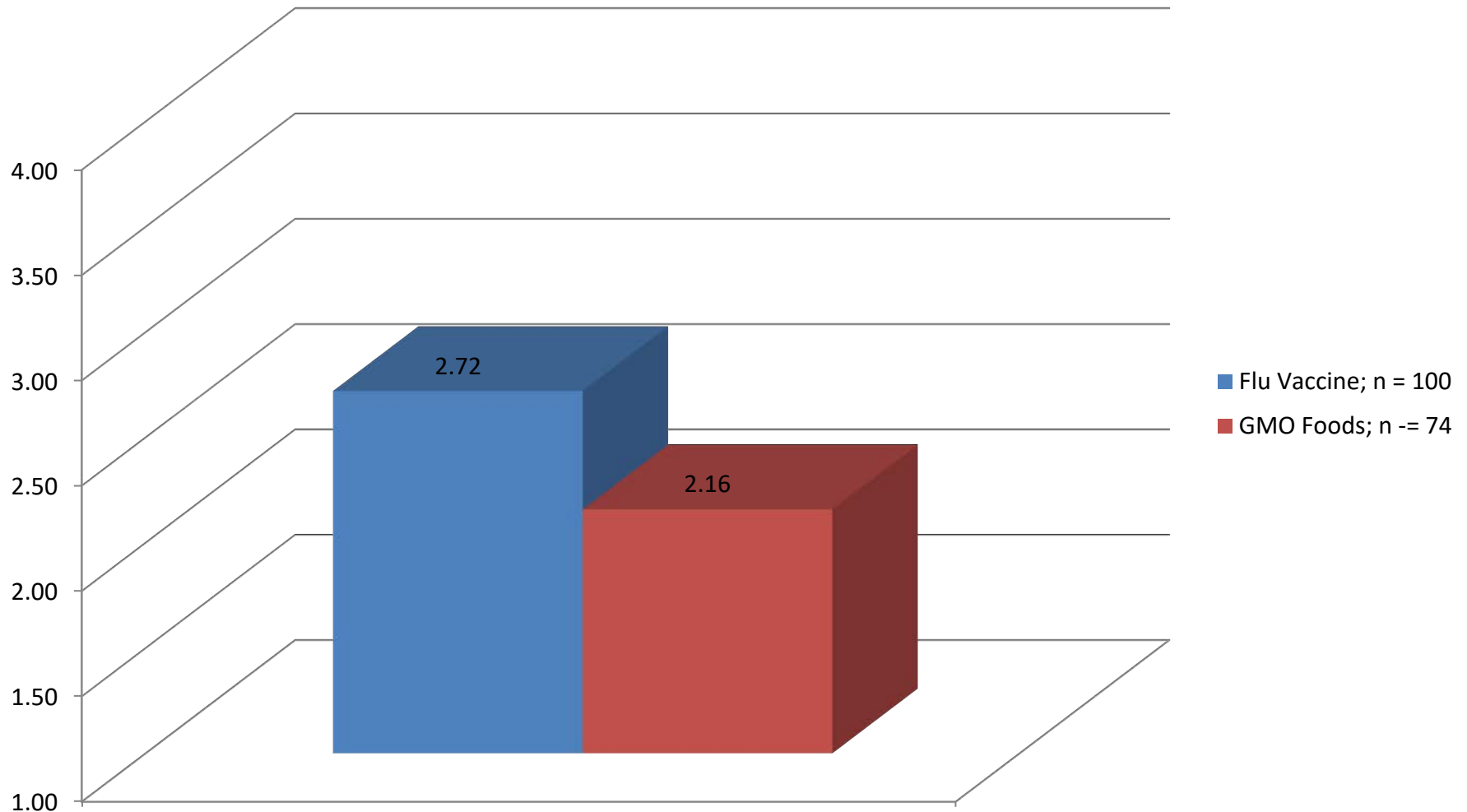
An independent samples t-test revealed no statistically significant differences between the means of these scenarios.



FYS Comparisons by Scenario for CT: Recommendation

Mean Scores on a scale of 1 – 4, with 4 being the highest possible score

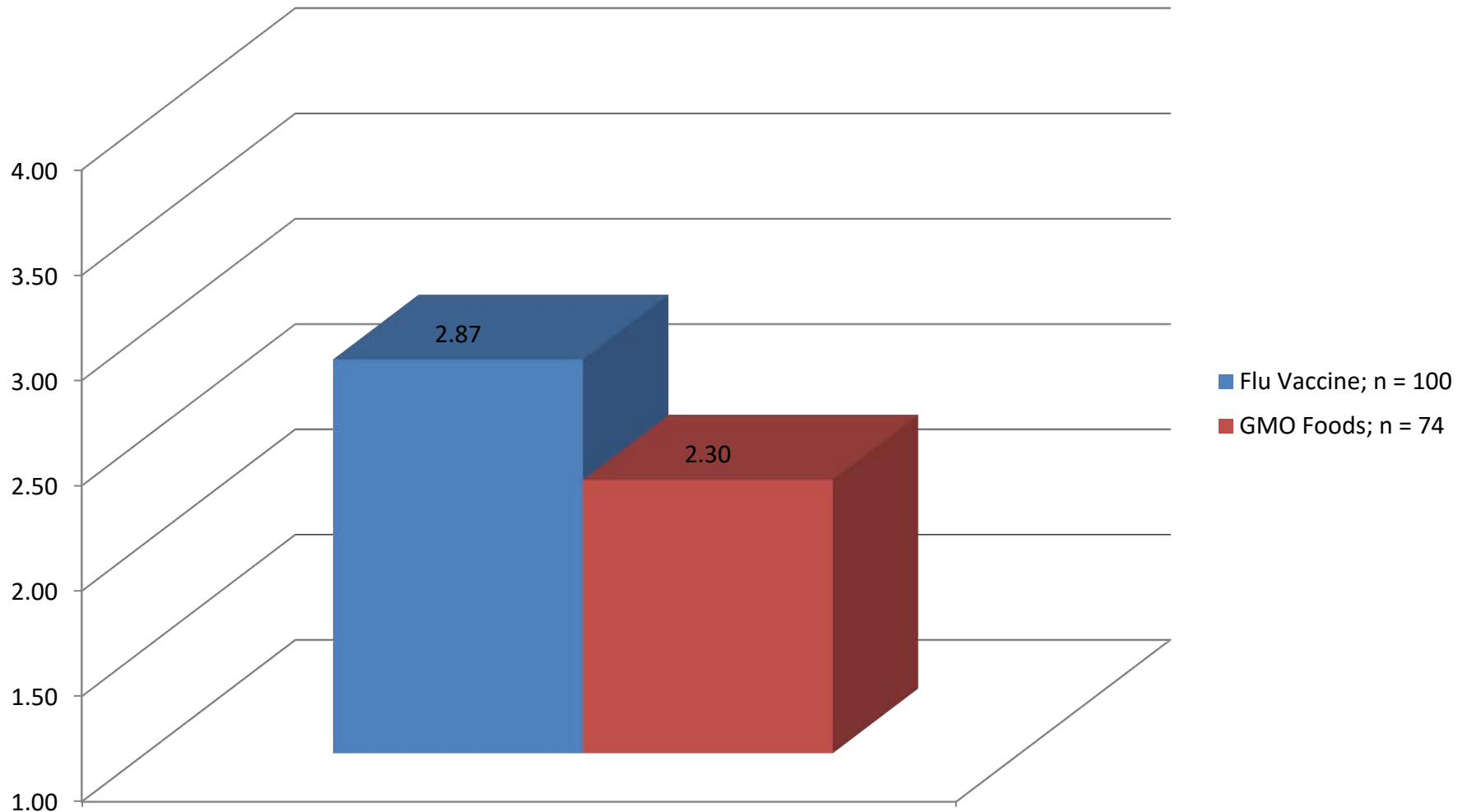
Using an adjusted alpha level of .017, an Independent Samples t-test revealed a statistically significant difference between the Flu Vaccine and GMO Foods Scenarios, $t(1429.960) = 5.210$, $p < .001$.



FYS Comparisons by Scenario for CF: Development

Mean Scores on a scale of 1 – 4, with 4 being the highest possible score

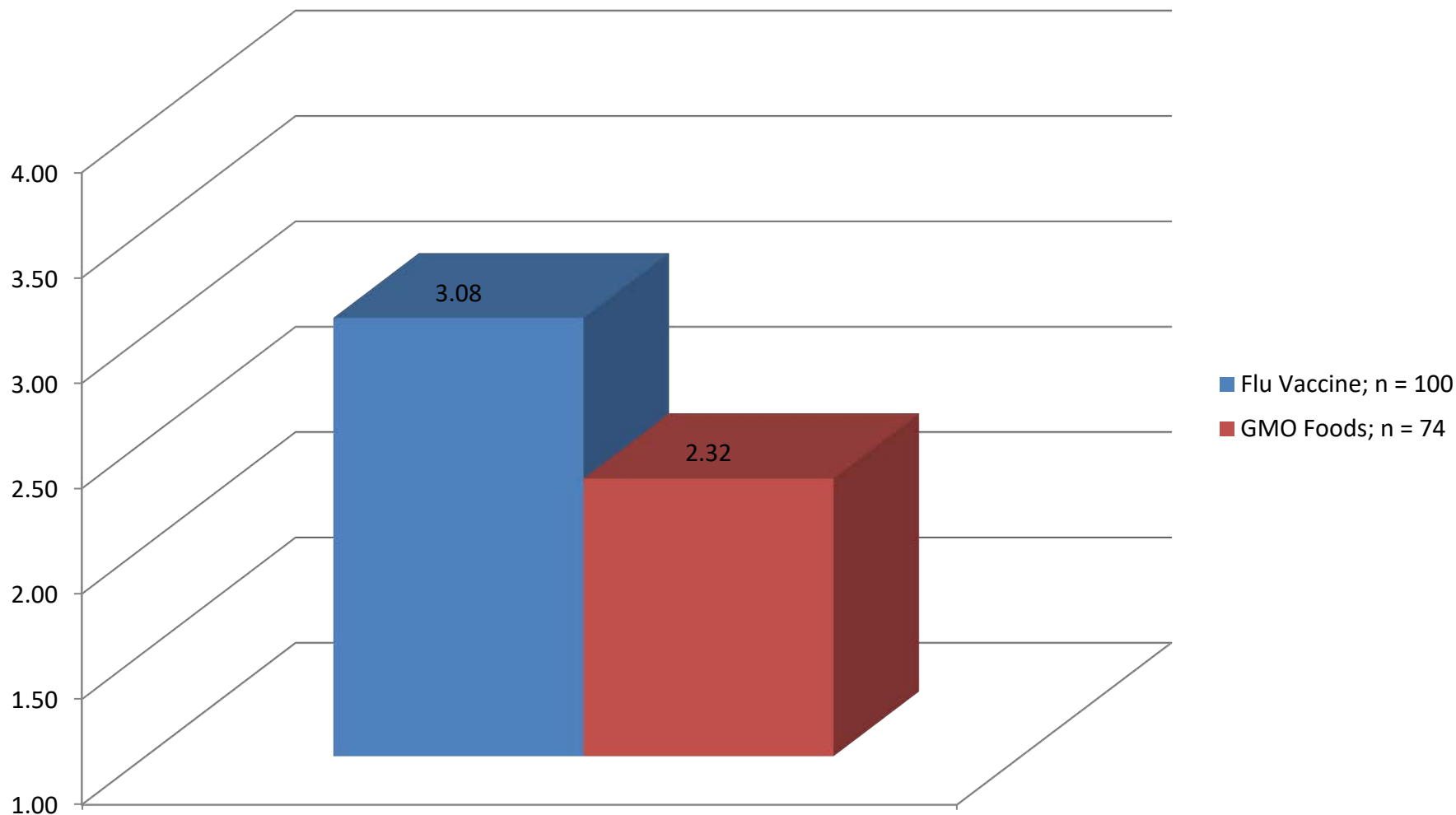
Using an adjusted alpha level of .017, an Independent Samples t-test revealed a statistically significant difference between the Flu Vaccine and GMO Foods Scenarios, $t(136.436) = 5.106, p < .001$.



FYS Comparisons by Scenario for CF: Convention/Format

Mean Scores on a scale of 1 – 4, with 4 being the highest possible score

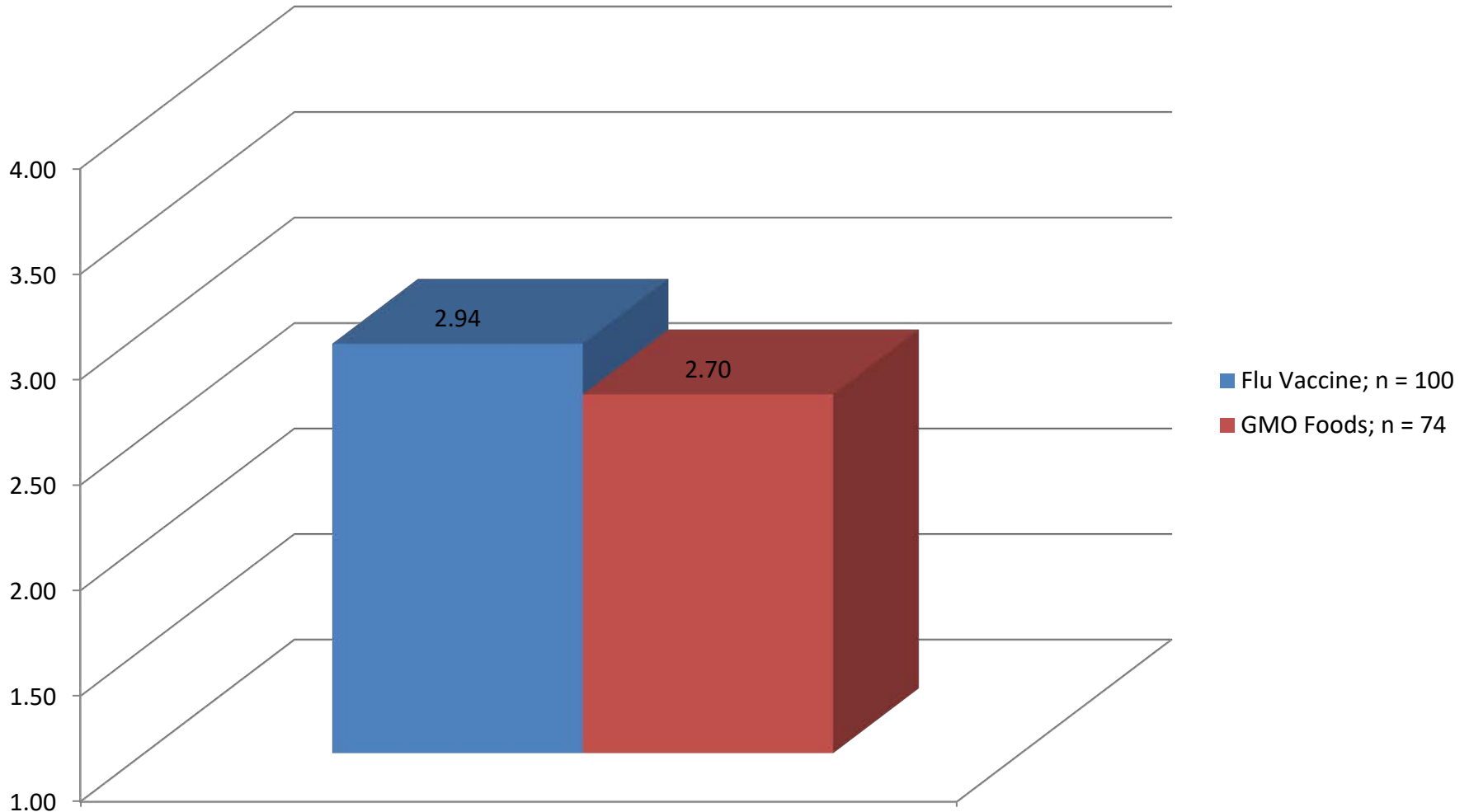
Using an adjusted alpha level of .017, an Independent Samples t-test revealed a statistically significant difference between the Flu Vaccine and GMO Foods Scenarios, $t(129.808) = 7.285, p < .001$.



FYS Comparisons by Scenario for CF: Communication Style

Mean Scores on a scale of 1 – 4, with 4 being the highest possible score

Using an adjusted alpha level of .017, an Independent Samples t-test revealed a statistically significant difference between the Flu Vaccine and GMO Foods Scenarios, $t(123.260) = 3.447, p < .001$.



Reference

Stellmack, M.A., Kohneim-Kalkstein, Y. L, Manor, J. E., Massey, A. R., & Schmitz, J. A. P. (2009). An assessment of reliability and validity of a rubric for grading APA-style introductions. *Teaching of Psychology*, 36, 102-107.