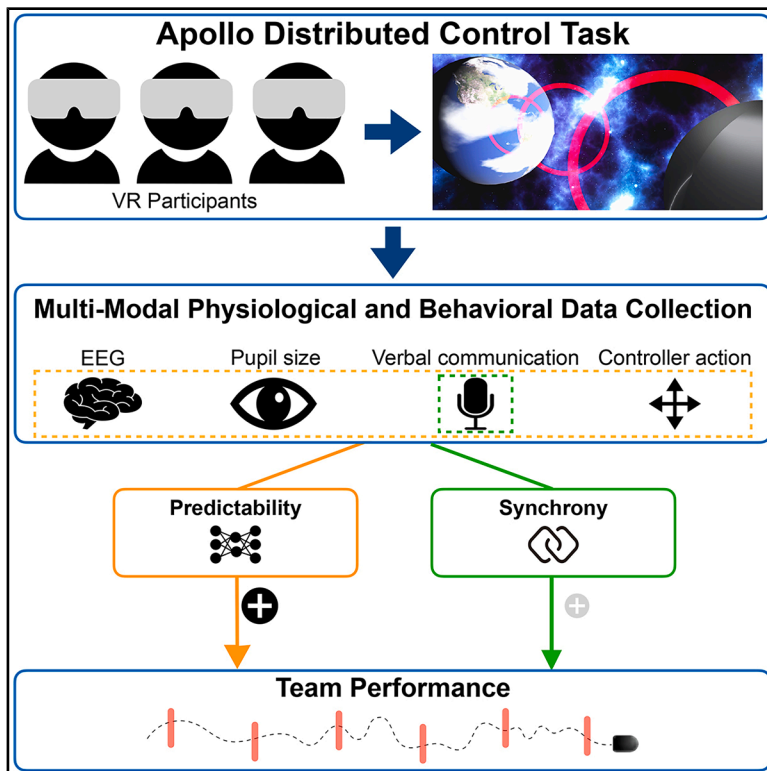


Physiologically informed predictability of a teammate's future actions forecasts team performance

Graphical abstract



Authors

Yinuo Qin, Richard T. Lee, Weijia Zhang, Xiaoxiao Sun, Paul Sajda

Correspondence

yinuo.qin@columbia.edu (Y.Q.),
psajda@columbia.edu (P.S.)

In brief

Artificial intelligence; Social sciences

Highlights

- Developed a VR-based sensorimotor task to study triadic human teaming
- Used multimodal deep learning to identify a biomarker linked to team performance
- Showed that this biomarker predicts team performance better than synchrony
- Findings provide insights into optimizing team dynamics in complex settings



Article

Physiologically informed predictability of a teammate's future actions forecasts team performance

Yinuo Qin,^{1,4,*} Richard T. Lee,^{1,2} Weijia Zhang,¹ Xiaoxiao Sun,¹ and Paul Sajda^{1,2,3,*}¹Department of Biomedical Engineering, Columbia University, New York, NY 10027, USA²Department of Electrical Engineering, Columbia University, New York, NY 10027, USA³Department of Radiology, Columbia University, New York, NY 10027, USA⁴Lead contact*Correspondence: yinuo.qin@columbia.edu (Y.Q.), psajda@columbia.edu (P.S.)<https://doi.org/10.1016/j.isci.2025.112429>

SUMMARY

In collaborative environments, a deep understanding of multi-human teaming dynamics is essential for optimizing performance. However, the relationship between individuals' behavioral and physiological markers and their combined influence on overall team performance remains poorly understood. To explore this, we designed a triadic human-collaborative sensorimotor task in virtual reality (VR) and introduced a predictability metric to examine team dynamics and performance. Our findings reveal a strong connection between team performance and the predictability of a team member's future actions based on other team members' behavioral and physiological data. Contrary to conventional wisdom that high-performing teams are highly synchronized, our results suggest that physiological and behavioral synchronizations among team members have a limited correlation with team performance. These insights provide a quantitative framework for understanding multi-human team dynamics, paving the way for deeper insights into team dynamics and performance.

INTRODUCTION

Teamwork is a critical form of human interaction, productivity, and survival. From world championship sports teams to intimate working groups, from ancient tribal rituals to modern urban planning, teaming has consistently been a critical and innate element of human behavior. Without teaming, our society would likely look very different, lack rich and diverse cultures, lack marvels of engineering and construction, and have limited groundbreaking scientific advancements. Studying the fundamental mechanisms behind human teaming is essential to understanding and improving collective human intelligence.

Games and collaborative tasks have been used as major platforms to study multi-human teaming. From role-playing to battle arena games, many previous studies have shown that multiplayer online games have great potential for studying team dynamics,¹ leadership in multi-human teaming,^{2,3} and individual behavior within teams.⁴ However, it is still unclear whether the insights gained from simple game-based studies can be generalized to more complex, high-stakes team interactions and team performance.

In addition to computer games, real-world scenarios, such as simulated hospitals with surgical teams and teams in manufacturing companies, have been used to study team performance and effectiveness.^{5,6} Most of these previous studies have used qualitative methods such as interviews,⁶ question-

naires,^{6,7} and surveys.^{5,8} While these qualitative studies can help us gain insights into how some task-related factors can impact team performance, these methods are prone to bias due to subjective reporting and are often difficult to reproduce. Therefore, the additional consideration of quantitative evaluation metrics to understand team performance remains essential but largely unexplored.

With the development of virtual reality (VR), more environmentally controlled team-based studies have been conducted.^{9–12} VR provides an immersive experience while reducing external distractions. The virtual environment also has the potential to provide realistic simulations with well-controlled delivery and simultaneous recording of events and interactions. However, few VR experiments have involved real-time synchronization and multi-modal data collection of multi-person teams. Most team-based VR experiments are conducted with a single human performing collaborative tasks with other simulated computer agents instead of working in the simulation with other people.^{10,11} These experiments limit the possibility of studying multi-human teaming.

Through studying human teaming in various tasks, previous research has highlighted that physiological synchrony among team members is positively correlated with team performance.^{13–16} Conversely, other studies have suggested a negative correlation between behavioral synchrony and team performance.^{17,18} The preceding literature lacks studies that



comprehensively correlate performance with both behavioral and physiological synchrony in complex teaming tasks. The interpretation of such correlations of team performance with behavioral synchrony and physiological synchrony remains unclear and incomplete. Therefore, we hypothesized that a comprehensive understanding of the balance between physiological and behavioral synchrony is critical for enhancing team performance, especially in tasks that demand high levels of cooperation, coordination, and collaboration.

In this work, we developed a framework to study multi-human teaming in a VR environment by quantitatively analyzing multi-modal physiological and behavioral data from all team members. We constructed an immersive sensorimotor task requiring three participants to collaboratively navigate a spacecraft, capturing multi-modal data from all participants (Videos S1, S2, S3, S4, S5, and S6). To identify potential biomarkers of team performance, we employed two computational approaches: inter-subject correlation (ISC) and predictability. ISC, traditionally linked to team performance metrics,^{19,20} was found to correlate with team performance only under specific measurements in our complex collaborative task. To address these limitations, we proposed a predictability approach, using a deep learning model to forecast one team member's remote controller actions based on their teammates' physiological and behavioral data. This predictive model revealed a significant correlation between the predictability of team members' actions and team performance, suggesting that predictability can serve as a robust biomarker for understanding and enhancing team dynamics in collaborative tasks.

RESULTS

Virtual reality paradigm for studying team dynamics

To test the correlation between team performance and physiological and behavioral synchrony among team members, we designed a multi-human team-based virtual reality (VR) task that we refer to as the **Apollo Distributed Control Task** (ADCT). Our task is inspired by the renowned Apollo 13 reentry mission and its extended cinematic story.^{21,22} The Apollo 13 mission is considered one of the history's most "successful failures" in that three astronauts exhibited extraordinary teamwork while operating different controls of a spacecraft collaboratively to safely navigate back to Earth after an oxygen tank exploded. The ADCT is a team-based version of a boundary avoidance task (BAT), which is a sensorimotor task requiring participants to navigate within a predefined area while avoiding boundary violations. Previous work has shown that BAT task induces arousal changes in individuals, affecting cognitive control and performance.²³ The ADCT has the following features built into its design and construction: 1) it is a challenging enough cooperative and collaborative task to trigger complex team dynamics; 2) the experiment was conducted repetitively with a consistent group of participants; 3) the task state and behavior of subjects are synchronized in real-time with simultaneously recorded multi-modal physiological signals; and 4) team performance is quantitatively assessed by evaluating the contributions of all team members, where local performance is measured in relation to short-term goals, and global perfor-

mance encompasses high-level planning tragedies. Figure 1 summarizes the ADCT.

Specifically, the ADCT is performed by a triad team in VR (Figure 1A). Each team member, as a co-pilot, has partial observation of the exterior space environment through uniquely positioned spacecraft windows, each with different viewing points. Each co-pilot controls a single degree of freedom of the spacecraft's movement, such as yaw, pitch, or thrust (Figure 1B). The team's goal is to collaboratively navigate the spacecraft back to Earth by following a predefined reentry path. The transparent red rings mark the boundary of the path, and the team must reach Earth within a limited time. Therefore, failing to pass all rings with sufficient speed results in trial failure. Teams are monetarily incentivized to complete as many trials successfully as possible. If they cannot return to Earth in time, they must navigate the entry path to get as close to Earth as possible.

While the teams performed the ADCT, we simultaneously collected electroencephalography (EEG), pupillometry, eye gaze, speech, and remote controller inputs from all participants (Figure 1C). Each team participated in three experimental sessions. The roles of participants were randomly assigned for each session, but the team members remained the same across all sessions. Each experimental session included 45 trials, each consisting of 15 rings. Team performance was quantitatively evaluated by the team's total number of ring obstacles successfully navigated.

Team performance improves across experimental sessions

We first analyzed performance dynamics across three experimental sessions to investigate how physiological and behavioral synchrony among team members relates to team performance. As shown in Figure 1D, the total number of rings passed by each team increased monotonically over the experimental sessions, indicating a steady improvement in overall team performance. Repeated measures analyses of variance (ANOVA) revealed significant performance differences across sessions ($F(2, 32) = 11.99, p < 0.001$). Post-hoc comparisons with Bonferroni correction showed a substantial improvement in performance from Session 1 to Session 3 ($p < 0.001$). Similarly, the averaged trial performance also improved significantly over time (Figure 1E, $F(2, 32) = 13.17, p < 0.01$). The performance significantly increased from Session 1 to Session 3 ($p < 0.001$). These findings suggest that team performance improved consistently as participants engaged in more task sessions. This steady enhancement highlights the potential for learning and adaptation in team dynamics through repeated collaborative tasks in immersive environments.

Subjective ratings and multi-modal inter-subject synchrony

After each experimental session, all co-pilots provided subjective ratings of their familiarity with and helpfulness toward other team members (see Post-Task Survey in [post task survey](#) for details). Analyzing these ratings allows us to track how familiarity and helpfulness change over time and investigate the potential impact of team members' perceptions on team performance. Surprisingly, the helpfulness rating shows a consistent decrease

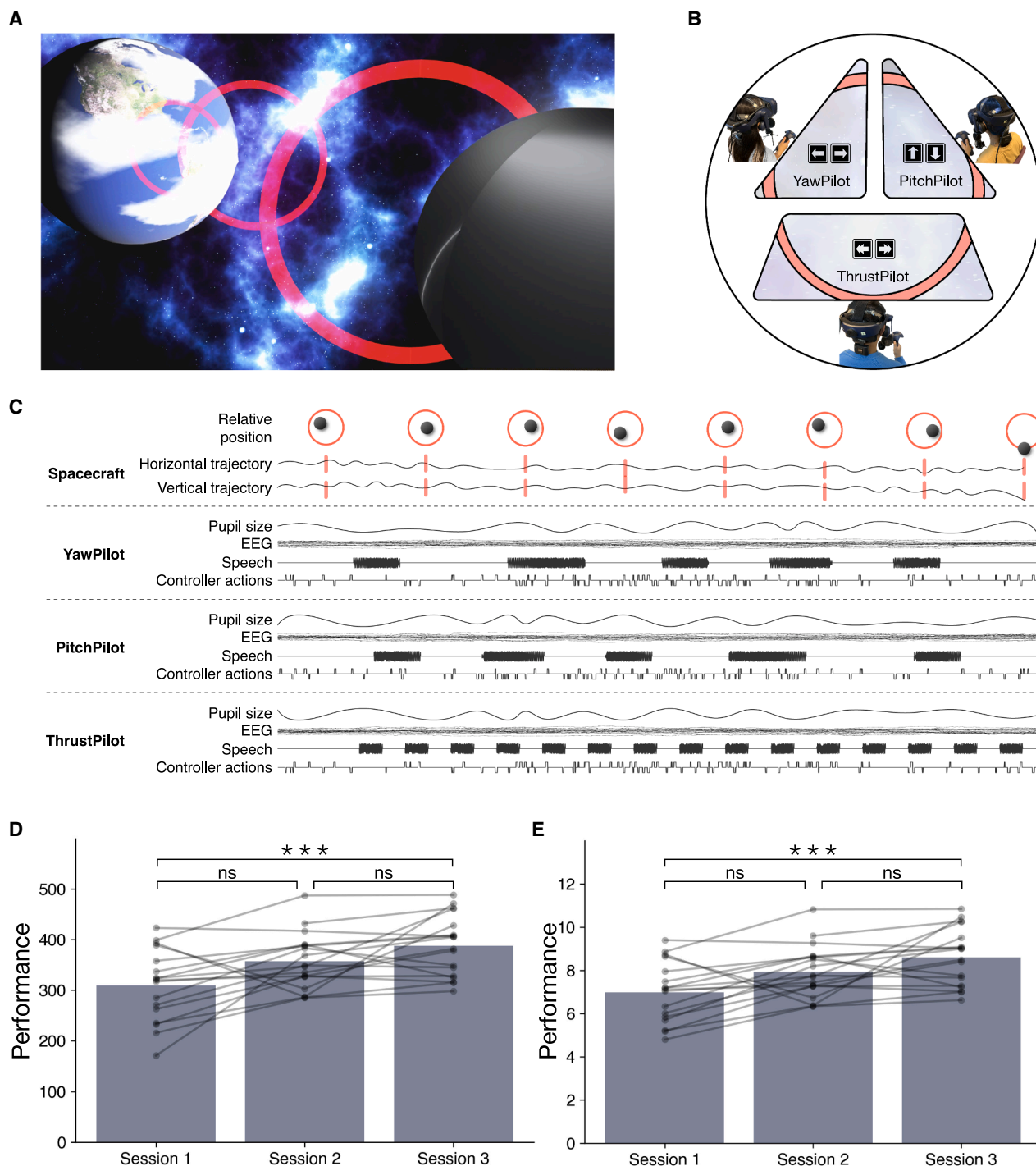


Figure 1. ADCT environment and team performance

(A) An illustration of the ADCT virtual environment. The team's goal is to control the spacecraft, passing all red rings and arriving back on Earth within the specified time limit.

(B) The view of three co-pilots with respect to a ring obstacle and the degree of freedom controlled by each role. The three co-pilots are YawPilot, PitchPilot, and ThrustPilot. Each participant was equipped with a VR headset, a microphone, a remote controller, and an EEG headset.

(C) Illustration of data modalities collected from all co-pilots. The red bars on the spacecraft's horizontal and vertical trajectories represent the relative location of ring obstacles. The uppermost section illustrates the cross-section of a spacecraft's position with respect to a ring.

(legend continued on next page)

across the experimental sessions (Figure 2A, repeated measures ANOVA, $F(2,34) = 9.33, p < 0.001$). In contrast to the decreasing helpfulness scores, the average familiarity rating across teams increases monotonically (Figure 2B, $F(2, 34) = 21.42, p < 0.001$). This pattern suggests that as team members become more familiar with each other, their perceptions of helpfulness may become more critical or nuanced.

Next, we analyzed the dynamics of team synchronization by calculating the inter-subject correlation (ISC) across various data modalities. ISC is a widely recognized metric for evaluating the synchrony among individuals performing identical tasks^{14,24–27} or collaborative tasks.^{28,29} This work analyzed the ISC among three co-pilots based on their pupil dynamics, EEG, remote controller inputs, and speech events. As illustrated in Figure 2C, pupil size synchrony remains relatively stable across sessions ($F(2,28) = 0.65, p = 0.53$). Interestingly, EEG ISC is maximized in the second experimental session. However, variations in EEG ISC did not achieve statistical significance (Figure 2D, $F(2, 18) = 1.51, p = 0.25$). Remote controller actions and speech events also remain stable along the three experimental sessions (Figures 2E and 2F; remote controller actions synchrony $F(2, 32) = 1.10, p = 0.34$; speech event synchrony $F(2,14) = 0.23, p = 0.80$). These findings suggest that increasing experimental sessions has a limited impact on synchronizations among team members' behavioral or physiological data.

Inter-subject synchrony and its correlation with team performance

Inter-subject synchrony (ISC) is often hypothesized to be correlated with team performance. Previous studies have demonstrated a positive relationship between team performance and synchrony in brain and pupil dynamics.^{14,19,20,30} However, whether synchrony among more than two team members correlates with overall team performance remains unexplored. To address this, we employed generalized linear mixed-effects models (GLMMs) to examine the relationship between inter-subject synchrony across multiple modalities and team performance (see [generalized linear mixed-effect model](#) for details).

Our findings reveal that behavioral synchrony, such as controller action synchrony and speech event synchrony, significantly correlates with team performance (Figure 2G). Interestingly, speech event synchrony among team members is positively correlated with team performance, suggesting that verbal communication enhances high-level task outcomes ($\beta = 1.63, P = 0.039$). In contrast, controller action synchrony is negatively correlated with team performance, possibly reflecting a detrimental effect of over-coordination on individual autonomy in control actions ($\beta = -1.01, P = 0.072$). Physiological synchrony, however, did not show a significant correlation with team performance (pupil size synchrony, $\beta = -0.73, P = 0.328$; EEG synchrony, $\beta = 0.11, P = 0.845$). These results show that behavioral synchrony is a

key predictor of team performance in triad teams, highlighting a previously overlooked factor in team performance research.

Quantifying team predictability using multi-modal physiological and behavioral data

A high-performing team consists of members who consistently engage in predictable interactions.³¹ This predictability results from a deep understanding and harmony within the team, making it easier for team members to anticipate one another's actions and reactions to each other. In this study, we used a multi-head attention model to quantitatively measure how the future actions of a teammate could be predicted from their teammates' physiology and behavior.

First, we epoched multi-modal physiological and behavioral data from 1.5 s before each ring-passing event (Figure 3A). The model received inputs from the initial 1 s of this epoch, where each input included the spacecraft's trajectory and the behavioral and physiological data of two co-pilots. The model's output was the generated prediction of the constructive 0.5-s controller action of the third co-pilot (0.5 s before passing the ring). On average, co-pilots made about 0.3 remote controller actions in that time period. We evaluated predictability at the team level by averaging the individual predictability scores across the three co-pilots. Since speech event synchrony significantly correlates with team performance ($P < 0.05$), we excluded speech event data from the model input to avoid potential bias. (The supplementary materials include results from a model incorporating speech input for comparison.) By analyzing team predictability, we demonstrate its potential as a biomarker significantly associated with overall team performance.

This model architecture is designed to integrate multi-modal data with varying temporal and spatial characteristics (Figure 3B). The inclusion of cross-modal attention layers ensures that dependencies between modalities are effectively captured, particularly when aligning trajectory information across diverse behavioral and physiological data sources. Additionally, self-attention layers within each modality help extract meaningful intra-modal patterns, such as EEG synchrony and pupil size dynamics. By combining the outputs from all modalities, the feedforward network synthesizes complementary features, creating a unified representation that encapsulates the interactions between physiological, behavioral, and environmental data.

The cross-attention mechanism further enhances this representation by linking the fused multi-modal features with the target modality, improving the accuracy of controller action predictions. This architecture leverages the unique contributions of each modality while ensuring robustness against noisy data. Moreover, its modular and adaptable design makes it well-suited for analyzing our complex multi-modal experiment.

We hypothesized that the predictability of team members' future controller actions would significantly correlate with team

(D and E) Team performance across three experimental sessions. The performance is measured by the number of rings passed. Each dot represents as mean of one team ($N = 17$). Bars indicate the average across teams. Asterisks indicate statistically significant differences, defined as $***P < 0.001$, while ns indicates the difference is not significant (repeated measures ANOVA with Bonferroni correction). (D) Total number of rings passed by each team in each session. (E) Averaged the number of rings passed by each team in each trial in three sessions.

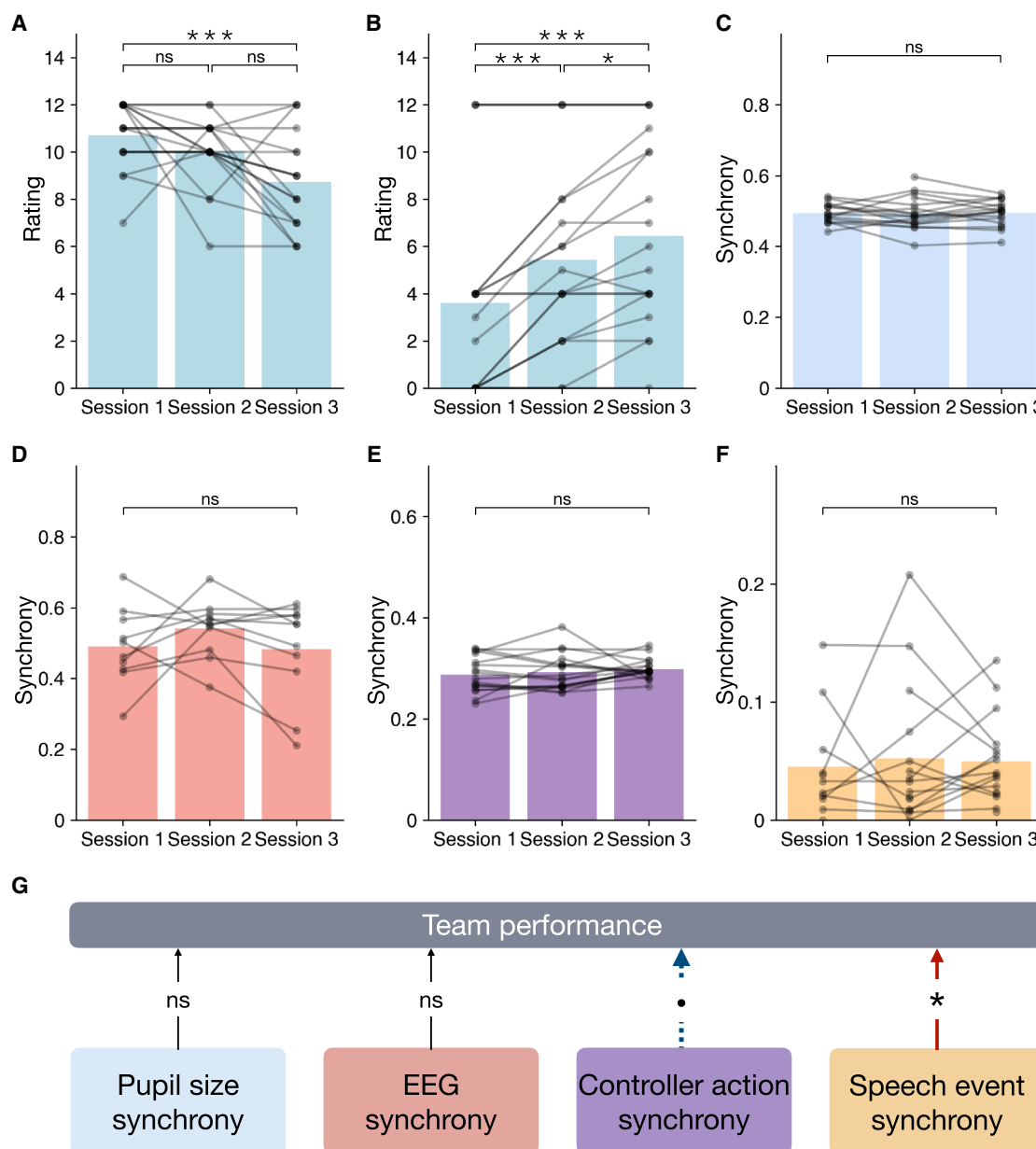


Figure 2. Subjective rating and multi-modal synchrony

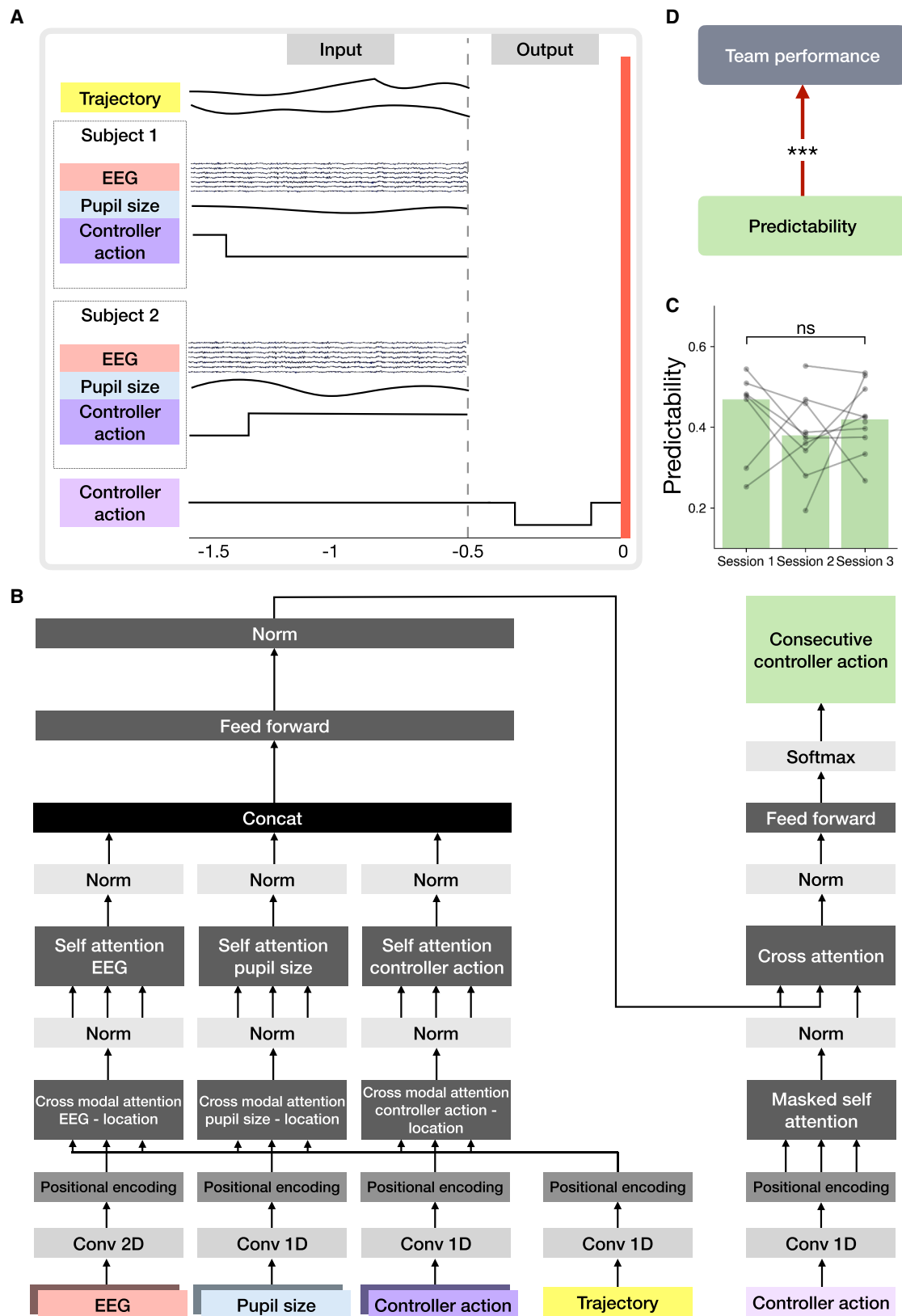
(A–F) Subjective ratings and synchrony among team members based on different physiological or behavioral data modalities across experimental sessions. Each dot is represented as a mean of one team, and the bars show the average across teams. Asterisks indicate statistically significant differences, defined as ns, not significant, $*P < 0.05$, $**P < 0.001$. Repeated measures ANOVA with Bonferroni correction. (A) Helpfulness rating of team members ($N = 17$). (B) Familiarity rating of team members ($N = 17$). (C) Pupil size synchrony among team members ($N = 14$). (D) EEG synchrony among team members ($N = 9$). (E) Remote controller action synchrony among team members ($N = 17$). (F) Speech event synchrony among team members ($N = 7$).

(G) Multi-modal synchrony and its correlation with team performance. Blue arrows indicate negative correlations, while red arrows indicate positive correlations. Asterisks indicate statistically significant differences, defined as ns, not significant, $*P < 0.1$, $**P < 0.05$. Generalized linear mixed-effects models.

performance. Consequently, we expected that the predictability of each team's actions would change across experimental sessions. As shown in Figure 3C, predictability changes slightly as the number of experimental sessions increases ($F(2, 10) = 0.12$, $P = 0.888$). A detailed analysis of predictability and team performance is provided in the next section.

Team action predictability as a performance biomarker

We identify an intriguing finding that predictability serves as a critical biomarker for team performance (Figure 3D, $\beta = 3.20$, $P < 0.001$). Specifically, a positive correlation between predictability and team performance suggests that when team members better anticipate each other's future actions, overall team



(legend on next page)

performance improves. This enhanced predictability may allow teams to execute more efficient joint strategies, leading to better outcomes, such as successfully passing more rings.

One possible explanation for this relationship is that increased predictability reduces uncertainty in decision-making, allowing team members to allocate cognitive and motor resources more effectively. In highly coordinated teams, members develop an implicit understanding of their partners' tendencies, minimizing reaction delays and facilitating smoother task execution.^{31,32}

Team performance and subjective ratings of team members

Familiarity among team members has been demonstrated to be positively correlated with team performance.^{15,33} Our experiment, focused on a collaborative distributed control task, observed a similar pattern (Figure 4). Specifically, familiarity among co-pilots is positively correlated with team performance, indicating that as familiarity increases, teams perform better on both sub-tasks and in achieving their long-term goals ($\beta = 0.18$, $P = 0.074$). This finding emphasizes the importance of building familiarity within teams, as it appears to enhance their ability to work cohesively and effectively toward shared goals.

Interestingly, the helpfulness rating of team members is significantly negatively correlated with team performance ($\beta = -0.32$, $P < 0.001$). This suggests that higher helpfulness ratings may reflect a greater reliance on teammates for support, which could reduce individual autonomy or efficiency, potentially detracting from the overall team performance. Conversely, lower helpfulness ratings may indicate a more balanced contribution from all members, optimizing team efficiency towards achieving shared objectives. Together, the subjective ratings of familiarity and helpfulness reveal a nuanced relationship between team dynamics and performance. While familiarity fosters cohesion and shared understanding, perceptions of helpfulness may introduce dynamics that negatively impact team performance. These insights highlight the complex interplay between subjective perceptions and performance, offering valuable guidance for designing and optimizing collaborative teams in distributed control tasks.

DISCUSSION

In this study, we conducted a team-based collaborative virtual reality (VR) experiment and demonstrated a multi-modal biomarker that directly correlates with team performance. Specifically, we demonstrated that a biomarker measuring the predictability of teammate behavior is better correlated with team performance. This biomarker is derived from integrating multi-

modal physiology and behavior of teammates to predict the future behavior of the remaining (i.e., left out) team member. Our predictability biomarker challenges the conventional wisdom that physiological and behavioral synchrony is a robust marker of a high-performing team.^{13,15–18}

Simultaneously collecting and analyzing multi-modal data is crucial for understanding team performance and dynamics. In contrast to executing simple tasks individually, collaborative tasks involve complex dynamics and interactions among team members. Various data modalities, including pupillometry, EEG, speech, and other physiological or behavioral data, have been analyzed individually but not in combination.^{34–37} We have developed a cross-modal multi-head attention predictive model that is capable of simultaneously analyzing multi-modal data from multiple team members (Figure 3B). This model integrates inputs from multiple data modalities, enabling not only the prediction of future actions of individuals but also the identification of a biomarker that is inversely related to overall team performance.³⁸ This result further demonstrates that different physiological and behavioral measures provide unique information that needs to be integrated to construct biomarkers that better relate to performance.

Our results revealed a positive correlation between our predictability biomarker and team performance. Aligning with the common belief that high-performance teams benefit from predictable actions among members,³¹ our findings suggest that this is expressed in teammate physiology in a way that leads to enhanced coordination and alignment for achieving higher performance (Figure 4). While the synchrony of individual modalities among co-pilots showed marginal or insignificant correlation with team performance, combining multi-modal data as input to the predictive model revealed that the predictability of team members' future actions is a stable and reliable biomarker of team performance. This highlights the potential of leveraging predictability as a key metric for understanding and improving team dynamics.

We have focused primarily on using predictability as a key indicator of team performance in collaborative tasks involving multiple humans. A pivotal question arising from our research is how we may practically leverage the predictive abilities of team members to enhance team dynamics and performance. This capability can facilitate collaboration between humans or, potentially, teaming between humans and artificial intelligence (AI) agents.^{39–41} Our findings lay the groundwork for innovative teaming strategies, fostering enriched and more productive collaborations. Another important direction for future work is to investigate how task difficulty levels impact team decision-making, as well as the physiological and behavioral responses of

Figure 3. Predictability of each team member's actions as a biomarker of team performance

(A) An illustration of a single epoch of the multi-modal data. Each epoch is relative to a ring, and we divided each epoch into input and output for the predictive model. The predicted action of an individual is based on a generative model that uses the behavioral and physiological data of the other two co-pilots. Predictability is evaluated by computing the correlation of the true action of a co-pilot with the model-predicted action.

(B) Multi-head attention modal structure. The cross-modal attention layers take the spacecraft trajectories and physiological or behavioral data.

(C) Team predictability across three experimental sessions. Each dot represents as mean of one team, and the bars show the average across teams ($n = 10$). ns, not significant (repeated measures ANOVA with Bonferroni correction.).

(D) Correlation between team performance and predictability. The red arrow indicates positive correlations. Asterisks indicate statistically significant differences, defined as $***P < 0.001$ (generalized linear mixed-effects models.).

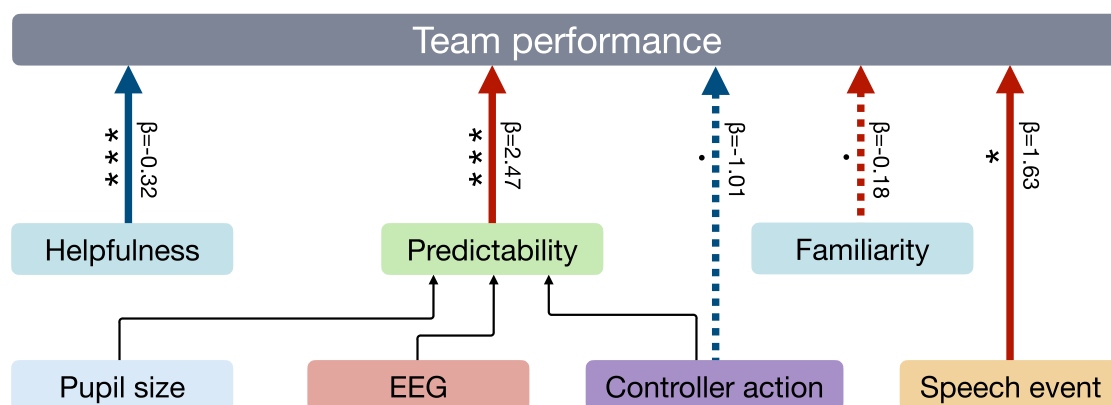


Figure 4. Overview of the correlation between predictability and team performance

All correlations are the GLMM results in accounting for session-level differences. The predictability biomarker is computed based on multi-modal physiological and behavioral data. Each arrow originates from the independent variable and points toward the dependent variable. The blue arrow indicates negative correlations, while the red arrow indicates positive correlations. Dashed lines indicate insignificant correlations. • $P < 0.1$, * $P < 0.05$, ** $P < 0.001$, $n = 10$, generalized linear mixed-effects models.

individuals.^{42,43} Understanding these effects could provide deeper insights into how teams adapt under varying cognitive demands and environmental constraints, shedding light on whether increased difficulty leads to heightened physiological synchrony, shifts in communication dynamics, or changes in collective decision-making strategies. Such analyses would further refine models of team coordination and improve the design of adaptive collaborative systems.

Limitations of the study

While our study provides valuable insights into team dynamics using a VR-based triadic sensorimotor task, several limitations should be noted. First, although our task offers an immersive and controlled setting, it differs from real-world scenarios in complexity and duration. Thus, further research should explore whether the findings generalize to longer, more complex collaborative tasks. Second, the predictability metric we introduced is based on short-term predictions (0.5-s future actions). Although this short-term prediction effectively captured immediate behavioral and physiological relationships among team members, its applicability to scenarios involving long-term strategic planning or more intricate decision-making remains unknown.

RESOURCE AVAILABILITY

Lead contact

Requests for further information and resources should be directed to and will be fulfilled by the lead contact, Yinyu Qin (yinyu.qin@columbia.edu).

Materials availability

This study did not generate new materials.

Data and code availability

- The preprocessed and de-identified data have been deposited at OSF: https://osf.io/u89s6/?view_only=8d7f6086c42f4bedb06232a12ac128b3 and are publicly available as of the date of publication. All other data reported in this article will be shared by the [lead contact](#) upon request.

- All original code has been deposited at https://github.com/liinc-lab/predictability_performance_and_ISC and is publicly available as of the date of publication.
- Any additional information required to reanalyze the data reported in this article is available from the [lead contact](#) upon request.

ACKNOWLEDGMENTS

This work was supported by funding from the Army Research Laboratory's STRONG Program (W911NF-19-2-0139, W911NF-19-2-0135, W911NF-21-2-0125), the National Science Foundation (IIS-1816363, OIA-1934968), the Air Force Office of Scientific Research (FA9550-22-1-0337), and a Vannevar Bush Faculty Fellowship from the US Department of Defense (N00014-20-1-2027).

AUTHOR CONTRIBUTIONS

Y.Q. and P.S. conceived the research program, designed experiments, and wrote the article; Y.Q. performed data collection and analysis with help from R.T.L., W.Z., X.S., and Y.Q. performed GLMM analysis; all authors made contributions to editing the article; P.S. supervised the entire work.

DECLARATION OF INTERESTS

The author Paul Sajda serves as a consultant for Optios Inc.

DECLARATION OF GENERATIVE AI AND AI-ASSISTED TECHNOLOGIES IN THE WRITING PROCESS

During the preparation of this work, the authors used ChatGPT to refine the text. After using this tool, the authors reviewed and edited the content as needed and take full responsibility for the content of the publication.

STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- [KEY RESOURCES TABLE](#)
- [EXPERIMENTAL MODEL AND SUBJECT DETAILS](#)
- [METHOD DETAILS](#)
 - Virtual environment

- Apparatus
- Procedure
- **QUANTIFICATION AND STATISTICAL ANALYSIS**
 - Data preprocessing
 - Post Task Survey
 - Pupil size, remote controller action, and speech event synchronies
 - EEG ISC
 - Statistical test
 - The generative forecasting model
 - Model evaluation
 - Predictability as a biomarker
 - Generalized linear mixed-effect model

SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.isci.2025.112429>.

Received: January 17, 2025

Revised: February 22, 2025

Accepted: April 10, 2025

Published: April 17, 2025

REFERENCES

1. Dabbish, L., Kraut, R., and Patton, J. (2012). Communication and commitment in an online game team. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pp. 879–888.
2. Jang, Y., and Ryu, S. (2011). Exploring game experiences and game leadership in massively multiplayer online role-playing games. *Br. J. Educ. Technol.* **42**, 616–623.
3. Pobiedina, N., Neidhardt, J., Moreno, M.d.C.C., Grad-Gyenge, L., and Werthner, H. (2013). On successful team formation: Statistical analysis of a multiplayer online game. In *2013 IEEE 15th Conference on Business Informatics (IEEE)*, pp. 55–62.
4. Sapienza, A., Zeng, Y., Bessi, A., Lerman, K., and Ferrara, E. (2018). Individual performance in team-based online games. *R. Soc. Open Sci.* **5**, 180329.
5. Morgan, P.J., Pittini, R., Regehr, G., Marrs, C., and Haley, M.F. (2007). Evaluating teamwork in a simulated obstetric environment. *Anesthesiology* **106**, 907–915.
6. Edmondson, A. (1999). Psychological safety and learning behavior in work teams. *Adm. Sci. Q.* **44**, 350–383.
7. Pearsall, M.J., and Ellis, A.P.J. (2006). The effects of critical team member assertiveness on team performance and satisfaction. *J. Manag.* **32**, 575–594.
8. Van der Vaart, T., and Van Donk, D.P. (2008). A critical review of survey-based research in supply chain integration. *Int. J. Prod. Econ.* **111**, 42–55.
9. Varlet, M., Filippeschi, A., Ben-Sadoun, G., Ratto, M., Marin, L., Ruffaldi, E., and Bardy, B.G. (2013). Virtual reality as a tool to learn interpersonal coordination: Example of team rowing. *Presence* **22**, 202–215.
10. Hansen, A., Larsen, K.B., Nielsen, H.H., Sokolov, M.K., and Kraus, M. (2020). Asymmetrical multiplayer versus single player: Effects on game experience in a virtual reality edutainment game. In *International Conference on Augmented Reality, Virtual Reality and Computer Graphics*, L.T. De Paolis and P. Bourdot, eds. (Springer), pp. 22–33.
11. Moore, S.M., and Geuss, M.N. (2020). Familiarity with teammate's attitudes improves team performance in virtual reality. *PLoS One* **15**, e0241011.
12. Weissker, T., and Froehlich, B. (2021). Group navigation for guided tours in distributed virtual environments. *IEEE Trans. Vis. Comput. Graph.* **27**, 2524–2534.
13. Henning, R.A., Boucsein, W., and Gil, M.C. (2001). Social-physiological compliance as a determinant of team performance. *Int. J. Psychophysiol.* **40**, 221–232.
14. Dikker, S., Wan, L., Davidesco, I., Kaggen, L., Oostrik, M., McClintock, J., Rowland, J., Michalareas, G., Van Bavel, J.J., Ding, M., and Poeppel, D. (2017). Brain-to-brain synchrony tracks real-world dynamic group interactions in the classroom. *Curr. Biol.* **27**, 1375–1380.
15. Gordon, I., Gilboa, A., Cohen, S., Milstein, N., Haimovich, N., Pinhasi, S., and Siegman, S. (2020). Physiological and behavioral synchrony predict group cohesion and performance. *Sci. Rep.* **10**, 8484.
16. Madsen, J., and Parra, L.C. (2022). Cognitive processing of a common stimulus synchronizes brains, hearts, and eyes. *PNAS nexus* **1**, pgac020.
17. Abney, D.H., Paxton, A., Dale, R., and Kello, C.T. (2015). Movement dynamics reflect a functional role for weak coupling and role structure in dyadic problem solving. *Cogn. Process.* **16**, 325–332.
18. Vicaria, I.M., and Dickens, L. (2016). Meta-analyses of the intra- and interpersonal outcomes of interpersonal coordination. *J. Nonverbal Behav.* **40**, 335–361.
19. Szymanski, C., Pesquita, A., Brennan, A.A., Perdakis, D., Enns, J.T., Brick, T.R., Müller, V., and Lindenberger, U. (2017). Teams on the same wavelength perform better: Inter-brain phase synchronization constitutes a neural substrate for social facilitation. *Neuroimage* **152**, 425–436.
20. Reiner, D.A., Dikker, S., and Van Bavel, J.J. (2021). Inter-brain synchrony in teams predicts collective performance. *Soc. Cogn. Affect. Neurosci.* **16**, 43–57.
21. Rerup, C. (2001). "Houston, We Have a Problem": Anticipation and Improvisation as Sources of Organizational Resilience (Snider Entrepreneurial Center, Wharton School).
22. (1995). *Apollo 13*. Universal Pictures Imagine Entertainment (Starring Tom Hanks, Kevin Bacon, and Bill Paxton).
23. Faller, J., Cummings, J., Saproo, S., and Sajda, P. (2019). Regulation of arousal via online neurofeedback improves human performance in a demanding sensory-motor task. *Proc. Natl. Acad. Sci. USA* **116**, 6482–6490. <https://doi.org/10.1073/pnas.1817207116>.
24. Kauppi, J.P., Jääskeläinen, I.P., Sams, M., and Tohka, J. (2010). Inter-subject correlation of brain hemodynamic responses during watching a movie: localization in space and frequency. *Front. Neuroinf.* **4**, 669.
25. Dmochowski, J.P., Sajda, P., Dias, J., and Parra, L.C. (2012). Correlated components of ongoing eeg point to emotionally laden attention—a possible marker of engagement? *Front. Hum. Neurosci.* **6**, 112.
26. Poulsen, A.T., Kamronn, S., Dmochowski, J., Parra, L.C., and Hansen, L.K. (2017). Eeg in the classroom: Synchronised neural recordings during video presentation. *Sci. Rep.* **7**, 43916.
27. Keles, U., Dubois, J., Le, K.J.M., Tyszka, J.M., Kahn, D.A., Reed, C.M., Chung, J.M., Mamelak, A.N., Adolphs, R., and Rutishauser, U. (2024). Multimodal single-neuron, intracranial eeg, and fmri brain responses during movie watching in human patients. *Sci. Data* **11**, 214.
28. Špiláková, B., Shaw, D.J., Czekóová, K., Mareček, R., and Brázdil, M. (2020). Getting into sync: Data-driven analyses reveal patterns of neural coupling that distinguish among different social exchanges. *Hum. Brain Mapp.* **41**, 1072–1083.
29. Xie, H., Karipidis, I.I., Howell, A., Schreier, M., Sheau, K.E., Manchanda, M. K., Ayub, R., Glover, G.H., Jung, M., Reiss, A.L., and Saggat, M. (2020). Finding the neural correlates of collaboration using a three-person fmri hyperscanning paradigm. *Proc. Natl. Acad. Sci. USA* **117**, 23066–23072.
30. Wohltjen, S., Toth, B., Boncz, A., and Wheatley, T. (2023). Synchrony to a beat predicts synchrony with other minds. *Sci. Rep.* **13**, 3591.
31. Bradley, B.H., Baur, J.E., Banford, C.G., and Postlethwaite, B.E. (2013). Team players and collective performance: How agreeableness affects team performance over time. *Small Group Res.* **44**, 680–711.
32. Butchibabu, A., Sparano-Huiban, C., Sonenberg, L., and Shah, J. (2016). Implicit coordination strategies for effective team communication. *Hum. Factors* **58**, 595–610.

33. Harrison, D.A., Mohammed, S., McGrath, J.E., Florey, A.T., and Vanderstoep, S.W. (2003). Time matters in team performance: Effects of member familiarity, entrainment, and task discontinuity on speed and quality. *Pers. Psychol.* 56, 633–669.
34. Stevens, R., Galloway, T., Gorman, J., Willemsen-Dunlap, A., and Halpin, D. (2016). Toward objective measures of team dynamics during healthcare simulation training. In *Proceedings of the International Symposium on Human Factors and Ergonomics in Health Care*, 5, pp. 50–54. <https://doi.org/10.1177/2327857916051010>.
35. McCraty, R. (2017). New frontiers in heart rate variability and social coherence research: techniques, technologies, and implications for improving group dynamics and outcomes. *Front. Public Health* 5, 267.
36. Dias, R.D., Zenati, M.A., Stevens, R., Gabany, J.M., and Yule, S.J. (2019). Physiological synchronization and entropy as measures of team cognitive load. *J. Biomed. Inform.* 96, 103250.
37. Dunbar, T.A., and Gorman, J.C. (2020). Using communication to modulate neural synchronization in teams. *Front. Hum. Neurosci.* 14, 332.
38. Madsen, A.G., Lehn-Schiøler, W.T., Jónsdóttir, Á., Arnardóttir, B., and Hansen, L.K. (2023). Concept-based explainability for an EEG transformer model. In *2023 IEEE 33rd International Workshop on Machine Learning for Signal Processing (MLSP) (IEEE)*, pp. 1–6.
39. Chen, B., Vondrick, C., and Lipson, H. (2021). Visual behavior modelling for robotic theory of mind. *Sci. Rep.* 11, 424.
40. Metcalfe, J.S., Perelman, B.S., Boothe, D.L., and McDowell, K. (2021). Systemic oversimplification limits the potential for human-ai partnership. *IEEE Access* 9, 70242–70260.
41. Harris-Watson, A.M., Larson, L.E., Lauharatanahirun, N., DeChurch, L.A., and Contractor, N.S. (2023). Social perception in human-ai teams: Warmth and competence predict receptivity to ai teammates. *Comput. Hum. Behav.* 145, 107765.
42. Dörner, D., and Funke, J. (2017). Complex problem solving: What it is and what it is not. *Front. Psychol.* 8, 1153.
43. Dörner, D., and Güss, C.D. (2022). Human error in complex problem solving and dynamic decision making: A taxonomy of 24 errors and a theory. *Computers in human behavior reports* 7, 100222.
44. Gramfort, A., Luessi, M., Larson, E., Engemann, D.A., Strohmeier, D., Brodbeck, C., Goj, R., Jas, M., Brooks, T., Parkkonen, L., and Hämäläinen, M. (2013). MEG and EEG data analysis with MNE-Python. *Front. Neurosci.* 7, 267. <https://doi.org/10.3389/fnins.2013.00267>.
45. Sainburg, T., Thielk, M., and Gentner, T.Q. (2020). Finding, visualizing, and quantifying latent structure across diverse animal vocal repertoires. *PLoS Comput. Biol.* 16, e1008228.
46. Makeig, S., Bell, A., Jung, T.P., and Sejnowski, T.J. (1995). Independent component analysis of electroencephalographic data. *Adv. Neural Inf. Process. Syst.* 8.
47. Winkler, I., Haufe, S., and Tangermann, M. (2011). Automatic classification of artifactual ica-components for artifact removal in eeg signals. *Behav. Brain Funct.* 7, 30.
48. Parra, L.C., Haufe, S., and Dmochowski, J.P. (2018). Correlated components analysis-extracting reliable dimensions in multivariate data. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1801.08881>.
49. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. *Adv. Neural Inf. Process. Syst.* 30.
50. Breslow, N.E., and Clayton, D.G. (1993). Approximate inference in generalized linear mixed models. *J. Am. Stat. Assoc.* 88, 9–25.
51. Skipper, S., and Josef, P. (2010). statsmodels: Econometric and statistical modeling with python. 9th Python in Science Conference.
52. Seabold, S., and Perktold, J. (2010). statsmodels: Econometric and statistical modeling with python. In 9th Python in Science Conference.

STAR★METHODS

KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Critical commercial assays		
Physiological and behavioral data	This paper	https://osf.io/u89s6/?view_only=8d7f6086c42f4bedb06232a12ac128b3
Software and algorithms		
Unreal Engine 4.25.4	Epic Games	https://www.unrealengine.com/en-US
LabRecorder (1.14.0)		https://github.com/labstreaminglayer/App-LabRecorder
Python MNE (1.6.1)	Gramfort et al. ⁴⁴	https://mne.tools/1.7/index.html
Python noisereducer	Sainburg et al. ⁴⁵	https://github.com/timsainb/noisereducer
Analytic code	This paper	https://github.com/liinc-lab/predictability_performance_and_ISC
Other		
VIVE Pro Eye	HTC	https://www.vive.com/sea/product/vive-pro-eye/overview/
B-Alert X24	Advanced Brain Monitoring	https://www.advancedbrainmonitoring.com/products/b-alert-mobile
RTX A6000 GPU	NVIDIA	https://www.nvidia.com/en-us/design-visualization/rtx-a6000/

EXPERIMENTAL MODEL AND SUBJECT DETAILS

54 healthy human participants (age = 23.67 ± 3.34 year (mean \pm standard deviation); 27 females, 27 males) voluntarily participated in the three experiments. These participants were divided into 18 triad teams, and each team participated in three sessions on different days. All participants were allocated to experimental groups based on their availability. Due to incomplete sessions, data from one team were omitted from the final dataset. Data from four teams were omitted from the pupil size analysis due to invalid pupil size recordings of one or more co-pilots. EEG data from nine teams were excluded from the analysis due to error-prone recordings identified during preprocessing. Similarly, speech event data from ten teams were excluded because the speech event detection algorithm failed to extract speech events from one or more participants within the team. No participants or teams dropped out of the experiment due to motion sickness or other symptoms related to virtual reality. All participants had normal or corrected-to-normal visual acuity and gave informed consent before participating in each experiment. Human subject protocols were approved by the Columbia University Institutional Review Board.

METHOD DETAILS

Virtual environment

The virtual environment was built using *Unreal Engine 4.25.4*. The four main reactive objects in the virtual environment were 1. a spacecraft, 2. a countdown timer, 3. the rings, and 4. the Earth. As shown in Figure 1B, three viewing windows with different shapes and at different positions were placed at the front of the spacecraft. Each subject in the triad team was assigned to look through one window, and the degree of freedom the subject controls was fixed per experiment session, corresponding to its respective window. The ThrustPilot, who controlled the speed of the spacecraft, had the largest unobstructed field of view, which was located at the bottom of the spacecraft. The YawPilot, who controlled the left-right spacecraft movement, was located at the top-left of the spacecraft. The PitchPilot had a viewing window on the top-right and controlled the up-down movement of the spacecraft. Because the positions and shapes of the windows were different, subjects with different roles had partial and biased views of the environment. The field of view of the virtual camera of each co-pilot is 80° in Unreal Engine.

A countdown timer bar was displayed at the bottom of each window to indicate the remaining time for each trial. Initially, the timer bar was completely black. As time elapsed, the black portion of the bar gradually decreased, revealing an increasing white segment. This white segment represented the time that had passed and was inversely proportionate to the black portion, which showed the remaining time. Each trial had a maximum duration of 55 s. The timer would automatically stop and reset if the team either successfully navigated through all the rings and approached Earth, or failed to pass through any ring during the trial. Despite this, with the default speed of the spacecraft, teams would require at least 60 s to pass through all rings in a trial, presenting a significant challenge and requirement for active participation and collaboration with the ThrustPilot.

The rings were transparent red toruses that represented a trial's reentry path. At the beginning of each trial, a sequence of fifteen rings was generated, spaced equally but positioned at varying horizontal and vertical coordinates. The distance between any two

adjacent rings was 50,000 units in *Unreal Engine*. The Earth was positioned at the end of the path with 50,000 units from the final ring. The trial ended when the spacecraft, operated by the team of participants, successfully navigated through all rings and stopped in front of Earth. Upon successful completion, the term “Successful” will be displayed on each participant’s head-mounted display (HMD) for 1 s. Subsequently, a new trial will automatically be started.

Apparatus

In all experiments, each participant was equipped with VIVE Pro Eye head-mounted displays (HTC Corporation; resolution: 1440×1600 pixels per eye; refresh rate: 90 Hz), and an EEG device with 20 electrodes was placed in accordance with the international 10–20 system (Advanced Brain Monitoring B-Alert X24; sample rate: 256 Hz). A USB microphone was set in front of each subject to enable communication between subjects, and Mumble (version 1.4.230) was running locally on each desktop. We used LabRecorder (version 1.14.0) to collect the multi-modal data. Each head-mounted display is connected to a desktop with an Intel Core i9 CPU and an NVIDIA RTX 2070 Super GPU. The three desktops were connected to a local, secure WiFi network with a 2.6 Gbps router using client-server network protocols to communicate. The server was another desktop with an Intel Core i9 CPU and an NVIDIA RTX 2080 Super GPU.

Procedure

In each experiment, three participants arrived at the lab and watched an instructional video before the first session. Following the setup of the EEG devices, participants were escorted to three separate EEG recording chambers designed to block sound and electrical noise. These chambers were additionally acoustically shielded with 2-inch thick soundproofing foam to prevent echoes and minimize noise interference. We assisted the participants in setting up head-mounted displays and remote controllers.

Individual eye calibration commenced once each participant was fully equipped and settled. The calibration was conducted using the VIVE Pro Eye system. Each experiment began with five pilot trials following eye calibration, allowing subjects to familiarize themselves with the task environment before the commencement of data collection. A trial was terminated when the team failed to pass a ring due to a crash or a miss or if the time limit was exceeded. After the pilot trials, participants were notified via headphones that the experiment had officially started.

Each team participated in three repeated sessions of the same experiment. Each session was spaced at least 24 h apart, and no participant had participated in nor had familiarity with the task before their first session. Within each experiment session, roles were randomly assigned to the subjects. After each experimental session, all participants were asked to complete a post-task questionnaire separately (see the post-task survey in [post task survey](#) for details).

QUANTIFICATION AND STATISTICAL ANALYSIS

Data preprocessing

We implemented different pre-processing methods for various data modalities. For pupil size data, we first detected and removed blinks and artifacts. Then, we applied linear interpolation and Z-scored the pupil size data. This was followed by averaging the pupil size between the left and right eyes and a fourth-order Butterworth lowpass filter to remove high-frequency noise.

The EEG pre-processing included filtering the raw EEG data using fourth-order Butterworth bandpass filters with bands 0.5 Hz–100 Hz (MNE 1.6.1⁴⁴). Manual bad channel rejection was conducted to remove error-prone channels in each recording. Then we performed Independent Component Analysis (ICA)⁴⁶ and used the Multiple Artifact Rejection Algorithm (MARA)⁴⁷ to separate and reject artifact components.

Remote controller actions were first down-sampled to 60 Hz. Next, values greater than 0.5 were assigned a value of 1, values less than –0.5 were assigned a value of –1, and values between –0.5 and 0.5 were assigned a value of 0.

The speech preprocessing involved three steps. First, we applied the noise reduction function⁴⁵ to the speech recordings from each subject to remove background noise. Next, we used a simple voice activity detection function to extract speech events. Finally, the speech events were down-sampled to 60 Hz. All data modalities were then epoched based on the relative time to the respective rings and saved for analysis.

Post Task Survey

After each experimental session, all participants were asked to complete a survey comprised of demographic and subjective questions. In this study, our analysis concentrated on two specific subjective questions.

- (1) How helpful was each of your teammates in reaching the final goal?
- (2) How well did you know each of your team members before today?

Each participant was required to select one of three possible answers for each question that concerned every other team member, excluding themselves. These answers were scaled as 0 = Not at all, 1 = A little, and 2 = Very well. The responses to the helpfulness (Question 1) and familiarity (Question 2) questions were assessed based on the team and the specific experiment session. The

helpfulness and familiarity scores ranged from 0 to 12 for each team. A score of 0 indicated that all three participants rated 'Not at all' for each of the other two team members. In contrast, a score of 12 indicated that every participant rated 'Very well' for their teammates.

Pupil size, remote controller action, and speech event synchronies

This study computed the inter-subject correlation (ISC) across the three subjects using their pupil sizes, remote controller actions, and speech events. For each experiment session, we computed the Pearson Correlation Coefficient (r) between each pair of participants, participant a and participant b , with their distinct roles within the same team, for one data modality at a time, using 1

$$r_{a,b} = \frac{\sum_{i=1}^n (a_i - \bar{a})(b_i - \bar{b})}{\sqrt{\sum_{i=1}^n (a_i - \bar{a})^2 \sum_{i=1}^n (b_i - \bar{b})^2}}. \quad (\text{Equation 1})$$

n was the length of one epoch of data. The team ISC in one epoch r_{epoch} was the averaged ISC across three co-pilots.

EEG ISC

To assess inter-brain synchronization, we computed ISC using Correlation Component Analysis (CorrCA).⁴⁸ This approach involved utilizing linear combinations of EEG channels or EEG signals with other data modalities as separate channels that maximize the ISC on the data obtained from subjects within the same team. In our study, we employed an improved version of CorrCA to compute the correlation between multiple subjects within the same team while performing a collaborative control task. The EEG signals of each subject contained 20 channels, and the approach finds a weight vector w that maximizes the Pearson Correlation between subjects in the team.

The weight vector w determines which linear combination of different channels provided the most significant correlation among team members. Given the EEG signals of the three subjects, denoted as X_1, X_2 , and X_3 , where $X_n \in \mathbb{R}^{D \times T}$ with D representing the number of channels and T representing the number of time steps in an epoch, the weight vector w could be computed by:

$$w = \operatorname{argmax}_w \left(\frac{w^T R_{12} w}{\sqrt{w^T R_{11} w} \sqrt{w^T R_{22} w}} \right); \quad (\text{Equation 2})$$

where $R_{ij} = \frac{1}{T} X_i X_j^T$

We defined the within subject covariance as $R_w = \sum_i^N R_{ij}$ and between subject covariance as $R_b = \sum_i^N \sum_{j>i}^N R_{ij}$. Here, $N = 3$ denoted the number of subjects in each experiment. We computed the eigenvectors e_k of $R_w^{-1} R_b$ and ranked the eigenvectors in descending order based on the corresponding eigenvalues. Hence, the ISC was the maximum value of the strengths of correlations C_k , where

$$C_k = \frac{e_k^T R_b e_k}{e_k^T R_w e_k}. \quad (\text{Equation 3})$$

Statistical test

All statistical analyses of team performance, subjective ratings of team members, synchrony measures, and team predictability were performed using repeated measures analysis of variance (repeated measures ANOVA) with Bonferroni correction. In all analyses, the average of a team's measurement during a given experimental session was treated as one data point. The repeated measures design accounts for within-team variability across sessions, allowing for more accurate estimation of session effects.

Significance levels were indicated in the figures using asterisks, which are defined in each figure legend. Specifically, p values were categorized as follows: $*p < 0.05$, $**p < 0.01$, and $***p < 0.001$. All statistical tests reported in the figures refer to the results of repeated measures ANOVA unless otherwise stated in the figure legends.

The generative forecasting model

The predictive model we implemented was a multi-head attention-based neural network that tracked relationships between events in data within the time domain. Figure 3B illustrates the structure of the model. The input to the model included the team's spacecraft trajectory along with the behavioral and physiological data of two participants. The transformer model utilized both encoders and decoders discussed in the original transformer model.⁴⁹ The 8-head attention layers in the encoder and the masked 8-head attention layers in the decoder were implemented as follows:

$$\begin{aligned} \text{MultiHead}(Q, K, V) &= \text{Concat}(\text{head}_1, \dots, \text{head}_4) W^O, \\ \text{where } \text{head}_n &= \text{Attention}(QW_n^Q, KW_n^K, VW_n^V), \\ [W_n^Q, W_n^K] &\in \mathbb{R}^{d_m \times d_k}, W_n^V \in \mathbb{R}^{d_m \times d_v}, W^O \in \mathbb{R}^{hd_v \times d_m} \end{aligned} \quad (\text{Equation 4})$$

We used $d_k = d_v = d_m/h = 64$ in this work. The *Attention* function took a set of queries as a matrix Q , the keys matrix K , and the values matrix V . The output of the *Attention* layer was:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V. \quad (\text{Equation 5})$$

All training and testing were conducted on a single NVIDIA RTX A6000 GPU, utilizing CUDA version V12.2.140. To further validate our model, we monitored metrics such as loss and accuracy during the training phase and utilized a validation dataset to assess performance periodically.

Model evaluation

We evaluated the predictive model's performance by computing the Pearson correlation coefficient r between the prediction and the target. To do so, we first computed the correlation of each individual using 1, where a was the concatenated target actions, and b was the concatenated model predictions.

Predictability as a biomarker

The predictive model we developed generates predicted future actions for each co-pilot based on the behavioral and physiological data of the other two co-pilots. These predictions are then correlated with the co-pilots' actual actions to compute a unique correlation score for each individual. We employ (1) to calculate the holistic team biomarker, which averages the predictability scores across the three co-pilots.

Generalized linear mixed-effect model

As an extension to the generalized linear model (GLM), the linear predictors of the generalized linear mixed-effects model (GLMM) contained random effects in addition to the usual fixed effects.⁵⁰ Given that our dependent variable—team performance—is measured as the count of successfully passed rings, it does not follow a normal distribution. Therefore, a GLMM was chosen.

We used the GLMM in Python statsmodels⁵¹ to investigate the relationship between varied variables with team difference considered as random-effect.⁵² The final regression formula of each model was listed in [supplemental information](#) in general form:

$$y = X\beta + Z\mu + \varepsilon, \quad (\text{Equation 6})$$

where y is the outcome variable. X represents the predictor variables. β is a column vector of the fixed-effects regression coefficients, and Z is the design matrix for the random effects (the random complement to the fixed X). μ is a vector of the random effects (the random complement to the fixed β), and ε is a column vector of the residuals.

Asterisks indicate statistically significant differences, defined as ns, not significant, $\cdot P < 0.1$, $*P < 0.05$, $**P < 0.001$.