

Meeting Effectiveness and Inclusiveness: Large-scale Measurement, Identification of Key Features, and Prediction in Real-world Remote Meetings

YASAMAN HOSSEINKASHI, Microsoft, USA
LEV TANKELEVITCH, Microsoft Research, United Kingdom
JAMIE POOL, Microsoft, USA
ROSS CUTLER, Microsoft, USA
CHINMAYA MADAN, Microsoft, USA

Workplace meetings are vital to organizational collaboration, yet relatively little progress has been made toward measuring meeting effectiveness and inclusiveness at scale. The recent rise in remote and hybrid meetings represents an opportunity to do so via computer-mediated communication (CMC) systems. Here, we share the results of an effective and inclusive meetings survey embedded within a CMC system in a diverse set of companies and organizations. We correlate the survey results with objective metrics available from the CMC system to identify the generalizable attributes that characterize perceived effectiveness and inclusiveness in meetings. Additionally, we explore a predictive model of meeting effectiveness and inclusiveness based solely on objective meeting attributes. Lastly, we show challenges and discuss solutions around the subjective measurement of meeting experiences. To our knowledge, this is the largest data-driven study conducted after the pandemic peak to measure, understand, and predict effectiveness and inclusiveness in real-world meetings at an organizational scale.

CCS Concepts: • **Human-centered computing** → **Empirical studies in collaborative and social computing**; • **Computer systems organization** → **Embedded systems**.

Additional Key Words and Phrases: Computer-mediated communication, meeting effectiveness, meeting inclusiveness, statistical modeling, machine learning, workplace meetings

ACM Reference Format:

Yasaman Hosseinkashi, Lev Tankelevitch, Jamie Pool, Ross Cutler, and Chinmaya Madan. 2024. Meeting Effectiveness and Inclusiveness: Large-scale Measurement, Identification of Key Features, and Prediction in Real-world Remote Meetings. *Proc. ACM Hum.-Comput. Interact.* 8, CSCW1, Article 93 (April 2024), 39 pages. <https://doi.org/10.1145/3637370>

1 INTRODUCTION

Workplace meetings are vital to collaboration and coordination in organizations [4]. Such meetings may range from recurring team stand-ups to information-sharing, planning, decision-making, and brainstorming meetings, among others [1, 75, 83]. They form a crucial way in which individuals and organizations collaborate and engage in sensemaking, ritual, and strategic change, as well as

Authors' addresses: Yasaman Hosseinkashi, Microsoft, 1 Microsoft Way, Redmond, USA, YAHOSSEI@microsoft.com; Lev Tankelevitch, Microsoft Research, 21 Station Road, Cambridge, United Kingdom, lev.tankelevitch@microsoft.com; Jamie Pool, Microsoft, 1 Microsoft Way, Redmond, USA, Jamie.Pool@microsoft.com; Ross Cutler, Microsoft, 1 Microsoft Way, Redmond, USA, Ross.Cutler@microsoft.com; Chinmaya Madan, Microsoft, 1 Microsoft Way, Redmond, USA, Chinmaya.Madan@microsoft.com.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

2573-0142/2024/4-ART93 \$15.00

<https://doi.org/10.1145/3637370>

experience stress [19]. Meetings, and how they are run, affect both productivity and employee well-being [74, 75]. Given their importance, understanding and evaluating the quality of meetings has been a focus of the small yet growing field of meeting science [4, 20, 22, 24, 43, 75]. Two key dimensions of workplace meetings are meeting *effectiveness*, defined here pragmatically as the attainment of business goals, and meeting *inclusiveness*, defined here as the extent to which participants feel they have an opportunity to contribute and all voices have equal weight [4, 20, 22, 24, 43]. Measuring and understanding meeting effectiveness and inclusiveness, and their drivers, is the prerequisite to improving workplace meetings through better practice guidelines and interventions.

The need to improve workplace meetings has become more imperative after the swift rise in remote and hybrid work, triggered by the COVID-19 pandemic, which has both increased the number of meetings that people have and shifted more meetings to all-remote or hybrid formats [12, 54]. This shift, enabled by video conferencing and other computer-mediated communication (CMC) systems, also provides an opportunity to understand and improve workplace meetings in a deep and scalable way using the rich data and automation afforded by CMC systems. Here, we use novel, large-scale data from real-world meetings in a leading CMC system to address the following research questions:

- **RQ1: How can we accurately measure meeting effectiveness and inclusiveness with a survey at scale and in the real-world context of meetings, as they occur during the workday?** Measuring meeting effectiveness and inclusiveness enables identifying key drivers and their interaction (see RQ2 below). Moreover, measurement that is large-scale (i.e., with sufficient statistical power) and real-world (i.e., occurring in the local context of meetings) provides metrics necessary to track changes in an organization's meeting culture, enabling individual organizations to improve their meetings by contextually understanding their own meeting practices and evaluating their own interventions. Critically, deploying surveys in organizations inevitably has implications for data quality, hence our focus here on accurate measurement.
- **RQ2: What are the common drivers of meeting effectiveness and inclusiveness across organizations, as measured objectively, and how do they interact?** Although prior work has examined drivers like agenda use and meeting promptness (discussed in Section 2), the role and interaction of drivers like length, meetings' recurring status, and attendee participation, remains unclear. Moreover, being conducted before this era of hybrid work, most prior research has not explored the role of relevant factors like audio and video participation, among others. Finally, prior research has not measured these drivers objectively and during real-world workplace meetings, and has not studied whether they generalize across organizations with different meeting experiences. These insights are necessary to develop interventions to improve workplace meetings.
- **RQ3: Can a single statistical model be generalized to predict meeting effectiveness and inclusiveness for individual meetings, across organizations and industries, without relying on regular survey measurement?** Predicting meeting effectiveness and inclusiveness passively (i.e., without the need for regular survey data collection) can empower more organizations with effective metrics for tracking and improving meeting experiences. However, the feasibility and accuracy of model-based metrics for effectiveness and inclusiveness requires verification with large-scale training data and precise testing, which has not been done before.

At the heart of our contribution is scaling up the subjective measurement of effectiveness and inclusiveness in remote meetings¹ by using integrated survey ratings collected within a CMC system across multiple organizations, and linking these measurements with objective meeting attributes available via the same system (e.g., video usage, meeting size, meeting duration, etc.). To do this, we build on and substantially advance the previously developed multivariate graph model of factors associated with meeting effectiveness and inclusiveness introduced by Cutler et al. [24]. The graph model introduced in [24] is a descriptive² model that can identify a network of statistically meaningful correlations between survey ratings and meeting attributes. [24] demonstrates an initial proof-of-concept for a meeting effectiveness and inclusiveness survey that is integrated into a CMC system (in-client survey). The in-client data collection in [24] was conducted within a subset of one organization for a short time during the early stage of COVID-19, before the new norms of remote and hybrid work settled. While findings confirmed the usefulness of this approach to measure meeting effectiveness and inclusiveness, they are based on a limited set of telemetry and a specific target population with unknown generalizability to other organizations or time periods. Here, we significantly extend this work in several ways:

- (1) We implement a scalable measurement method in a CMC system that uses surveys to collect subjective ratings of meeting effectiveness and inclusiveness at the end of meetings that can be scaled to collect data across organizations of any size. Using this system, we collect a real-world meeting dataset from five large organizations across a range of industries. This dataset was collected in 2022, after the peak of the COVID-19 pandemic, when the workforce had time to adapt its behavior to the new hybrid work context and related changes in collaboration norms. The final dataset contains 15K ratings from the rollout of this in-client survey as a new feature of the CMC system.
- (2) Using our large-scale dataset, we leverage and expand our descriptive graph model, including the range of considered meeting attributes, to conduct analyses of telemetry-captured meeting attributes and examine their interactions to contextualize and refine their relationship with meeting participation, effectiveness, and inclusiveness.
- (3) We test whether our approach—the descriptive graph model based on integrated survey measurement and meeting telemetry—generalizes across organizations in different industries and with different sizes, teams, and therefore meeting experiences.
- (4) To ultimately obviate the need for meeting rating surveys and thereby expand the scalability of our approach, we explore a predictive model that can predict meeting effectiveness and inclusiveness based solely on telemetry-captured meeting attributes, using survey ratings as “ground truth” training data.
- (5) Given the central role of survey measurement in our approach, we report on and address a key data quality issue with subjective measurement of meeting effectiveness and inclusiveness: rating skew (i.e., the tendency to avoid poor ratings). We share experiments and analysis results that aim to address this significant data quality challenge.

Our descriptive modeling shows that, although the meeting rating baseline varies by organization, there exists a robust and consistent set of dependencies and priorities with a meaningful correlation with the effectiveness and inclusiveness of remote meetings across different organizations. This implies that the main factors and priorities related to meeting effectiveness and inclusiveness do

¹Remote meeting in this work refers to a meeting where participants join the meeting via a CMC system. Our data does not include telemetry about the location where participants join meetings (e.g., from home or a meeting room).

²Descriptive models refer to statistical models that are developed for the purpose of inferring relations among independent and outcome variables. On the other hand, predictive models are statistical models focused on inferring the values of outcome variables given the independent variables whenever the outcome variable is not known.

not strongly depend on the industry or other organization-specific factors. Specifically, we find that the strongest predictor of meeting inclusiveness and hence effectiveness is whether the attendees vocally participate in the conversation. We also find that turning on video in small meetings (less than 8 participants) is correlated with a 6% increase in the probability of participation. In the absence of video, using a headset corresponds to a 20% increase in the odds of participation. We also show that it is harder to maintain inclusive and effective meetings with a large number of attendees: every 2 new participants corresponds to a 1 percentage point drop (absolute) in the meeting effectiveness and inclusiveness rating. Lastly, while call reliability is necessary for enabling participation and an inclusive environment, it needs to be accompanied by call quality to ensure both inclusive and effective meetings.

Moreover, we demonstrate that the descriptive model readily generalizes across organizations in different industries, with different sizes, and thereby with different meeting experiences. This is crucial because it tests the generalizability not only of our meeting-related insights but also of our measurement and modeling approach, suggesting that it can be deployed to support a wide range of organizations. In contrast, we find that predictive models that try to predict the experience of a specific attendee have lower performance than non-specific attendee models. We also show that transferring the predictive model that estimates individual ratings from one organization to another can come with a considerable drop in accuracy.

Finally, given the central role of survey measurement in our approach, we show that survey rating skew poses a significant challenge in measuring meeting effectiveness and inclusiveness. This impacts the reliability and utility of metrics that companies can rely on to improve meeting culture. We demonstrate via analyses and experiments that rating skew can be induced by certain characteristics of the survey and its deployment, is common across organizations, and can be mitigated with survey design choices.

Our results shed light on key factors that should be considered when planning to improve meeting experiences and culture in large organizations that use CMC systems in the new era of hybrid work. We also demonstrate the value and feasibility of our measurement and modeling approach deployed in real-world contexts at an organizational scale.

The rest of this paper is organized as follows: Section 2 reviews the related work in measuring meeting effectiveness and inclusiveness in meetings joined with a CMC system. Section 3 describes the design and execution of the in-client survey for five different corporations from different industries, the dataset variables, and the modeling methodology. The model developed and data-driven insights about the characteristics of effective and inclusive meetings (EIM) are discussed in Section 4. Section 5 discusses the survey skew and its solutions. Lastly, we summarize our findings, discuss theoretical and design implications, and explore current limitations and future work.

2 RELATED WORK

2.1 Meeting design and meeting effectiveness

Studying meetings has been a key focus of the small but growing field of meeting science (recent reviews are given in [4, 55]). One survey-based method to understand the drivers of meeting effectiveness is to look for associations between meeting characteristics and participants' perceived meeting effectiveness for their recent or typical meetings. Leach et al. [43] show that, among other features, the use of an agenda, quality of facilities, and ending on time were correlated with perceived effectiveness, whereas meeting size and length were not (N = 958 survey). Attendee involvement in the meeting mediated the observed relationships. Cohen et al. [20] examined 'meeting quality', a construct similar to effectiveness (N=367 survey), finding that the top significant drivers are meeting space, size, starting promptness, lighting quality, and organization type. Although meeting

size was a significant factor here, length was again not. Allen et al. [5] show meeting size to be negatively correlated with perceived meeting effectiveness. Likewise, Standaert et al. [83] found meeting size and duration to be negatively correlated with perceived meeting effectiveness, though only for certain meetings (e.g., those using telepresence).

In summary, prior research has illuminated various key drivers of meeting effectiveness, though the role and interaction of meeting size and length remains unclear. More importantly, the role of participation in meetings and its interaction with other factors has not been rigorously studied [4]. Leach et al. [43] found that attendee involvement is a mediator of effectiveness, but this was measured using a survey asking about overall participation across attendees (e.g., “Participation is widespread among meeting attendees”), rather than measured objectively for individuals. This is important given that participation can vary widely among attendees, and is, therefore, difficult for one observer to retrospectively estimate overall participation [24], [47]. Additionally, no research has examined how the role of various meeting characteristics may differ between recurring and one-off meetings, given their typically different goals and structure [58].

Moreover, almost all of the above research has relied on surveys administered outside the real-world context of meetings. That is, participants are asked to recall their last meeting hours later or their ‘typical’ meeting experiences (e.g., [20, 43]). This can introduce biases related to information recency or availability. Using surveys for all measurements can also introduce common-method bias [67]. Indeed, meeting science researchers have called for moving beyond solely survey methodology [55]. However, objective measurement of meeting characteristics (e.g., using telemetry or sensors) remains rare. Constantinides et al. [23] showed that room pleasantness, as measured via sensor readings of light, temperature, etc., is correlated with meeting effectiveness. Cutler et al. [24] analyzed meeting telemetry and survey rating data to show that meeting inclusiveness, the comfortableness of participating, audio/video quality, and using screen sharing and video were correlated with meeting effectiveness. A graph model was constructed with an area under the curve (AUC) of 0.68, which is promising but suggests there are many more factors to include to better understand meeting effectiveness. Lastly, given the difficulty of collecting detailed survey data, almost all prior research has relied on small samples which limits the statistical power of analyses. With the exception of [24], to our knowledge, there have been no attempts to deploy such measurement at an organizational scale.

Equally importantly, the vast majority of prior research was conducted before the COVID-19 pandemic when norms shifted dramatically towards remote and hybrid work, which relies heavily on audio and video conferencing [12, 54]. However, prior research has focused on characteristics of meeting facilities, such as lighting and space [20, 23], rather than factors such as the use of headsets, and the choice of audio and video participation, which have become increasingly more relevant. We are not aware of studies that have examined the role of headsets and screen-sharing in workplace meetings, despite these being important factors affecting communication [73].

Unlike in-person meetings, remote and hybrid meetings enable participants to dynamically choose between audio and video participation throughout meetings. Studies showing the significance of video over audio conferencing are remarkably sparse and mixed. Veinott et al. [89] showed that video helps non-native speakers better negotiate than audio-only conferencing. However, Habash [36] showed that video added little or no additional benefit over audio-only conferencing for group perception and satisfaction in distributed meetings. Instead of measuring a task metric or satisfaction, Daly-Jones et al. [25] showed that video does improve conversational fluency and interpersonal awareness over audio-only meetings. Similarly, Tang et al. [87] showed that video conferencing system usage drops significantly when the video feature is removed, and that video is used to help mediate participants’ interaction and convey non-verbal information. Sellen [80]

showed when remote participants join an audio or video conference with a conference room, in-room participants produced more interruptions and fewer formal handovers of the floor than remote participants (i.e., reflecting a more natural flow of conversation). However, video did not improve the interruption or handover rate for remote participants compared to audio-only. Standaert et al. [83] showed that telepresence systems improved meeting effectiveness over audio and video conferencing systems, though they did not compare audio and video conferencing directly. More recently, Cutler et al. [24] showed that video usage is correlated with meeting effectiveness and inclusiveness. Research is needed to understand how various forms of remote participation interact with meeting design characteristics.

In summary, there is a need for research on meeting design and effectiveness that relies on objective metrics where possible, and that is large-scale, situated in the real-world context of workplace meetings, and is up-to-date for the new era of remote and hybrid work that depends on audio and video conferencing.

2.2 Measuring meeting effectiveness

While there has been significant research in better *understanding* meeting effectiveness, there is no common consensus on how to *measure* it. For example, [32] measures meeting effectiveness by the percentage of agenda tasks that are completed. [83] provided 19 business meeting objectives and used a survey with a 5-point scale (1: Not at all effective to 5: Very effective) on how different meeting modalities achieved the business meeting objectives. [59] measure meeting effectiveness using two items: goal attainment and decision satisfaction. [74] measure it for meetings in a typical week using a 6-item survey (5-point scale), including: “achieving your own work goals”, “achieving colleagues’ work goals”, and “promoting commitment to what was said and done in the meeting”, among others. [43] measure it for meetings in a typical week with a 3-item survey (5-point scale): “achieving your own work goals”, “achieving your colleagues’ goals” and “achieving your department’s / section’s / unit’s goals.”

[23] measure meeting effectiveness (termed “execution” in their work) with one survey question (7-point scale): “Did the meeting have a clear purpose and structure, and did it result in a list of actionable points?”. [24] used a single question for a large-scale email survey: “How effective was the meeting at achieving the business goals?” (1: “Very ineffective” to 5: “Very effective”). The survey was also integrated into a CMC system, albeit using a star rating for response options (hovering over each star highlighted the associated description of each option).

In summary, whereas most prior research has included multiple aspects in measuring effectiveness, such as the perspectives of multiple attendees (e.g., one’s own and colleagues’ work goals [74]), or the presence of multiple features (e.g., a clear meeting structure and actionable points [23]), other work [24, 83] has focused on a single, pragmatic aspect: the achievement of business goals. With the exception of [24], no prior research has deployed any measurement of meeting effectiveness at scale in a real-world meeting context where respondents may be busy or particularly influenced by work-related interactions. Further research is, therefore, necessary to understand the quality of meeting effectiveness ratings in such a context.

2.3 Measuring meeting inclusiveness

Meeting inclusiveness, the extent to which participants feel they have an opportunity to contribute and all voices have equal weight, is a key aspect that contributes to effectiveness [24], [21]. [56] reviews studies on how trust and member inclusion are factors that foster collaboration in teams, although not meetings specifically. There are many guides on how to have inclusive meetings, e.g., [66], though remarkably few studies that we are aware of actually measure inclusiveness.

More broadly, inclusiveness has been extensively studied for organizations. [7] used a large-scale (N=10,976) employee survey to build a structural equation model that shows how transformation leadership and diversity management correlate to an inclusive organizational culture (with inclusiveness measured using a six-question survey). This echoes the approach defined in [30] (Chapter 1). [65] defined a three-question survey on measuring team inclusion. [72] further studied the relationships between organizational and supervisory inclusiveness, citizenship behavior, and affective commitment. In [30] (Chapter 17), Lukensmeyer et al. provide a list of important characteristics of truly inclusive meetings (discussed in detail in [24], Section 3.1.2).

[23] measured *psychological safety* (originally defined in [28] as “the absence of interpersonal fear that allows people to speak up with work-relevant content”) using a survey question asking “Did you feel listened to during the meeting, or motivated to be involved in it?” (7-point scale). Inclusiveness is also related to “group process losses”, a set of interaction dynamics identified in collaboration engineering research [60]. These include concepts such as “evaluation apprehension” (measured using survey items such as “I felt apprehensive about expressing my ideas and findings to the rest of the group”), and “domination” (measured using survey items such as “I felt that there was at least one person in the group who tended to participate much more than the other team members”) [53, 60]. The collective term “group process losses” captures the idea that these dynamics, i.e., a lack of inclusiveness, impair meeting effectiveness [60].

[24] used a single question in a large-scale email survey: “How inclusive was the meeting? In an inclusive meeting, everyone gets a chance to contribute and all voices have equal weight” (1: “Not at all inclusive” to 5: “Very inclusive”). Another survey question was integrated into a CMC system: “Did you feel included in the meeting?” (a star rating, with hover text, 1: “I didn’t feel included at all” to 5: “I felt very included”).

The experience of inclusiveness may be related to individual factors such as gender. While there are many studies showing gender bias in speaking and interruption rates [29, 39, 44], both [20, 43] show that gender is not correlated to meeting effectiveness. [88] showed women felt more included and participated more when CMC meetings were used before face-to-face meetings, compared to after. [35] showed that traditional face-to-face meetings outperformed videoconferencing when accounting for team-building experience (but a dialogue-based framework in virtual teams can mitigate these differences).

For meeting effectiveness, with the exception of [24], no prior research has deployed any measurement of meeting inclusiveness at scale in a real-world meeting context. Further research is therefore needed to understand the quality of meeting inclusiveness ratings in such a context, particularly given its relatively more sensitive nature [21].

2.4 Predicting meeting effectiveness and inclusiveness

Given the challenges of collecting large-scale survey data, using passively and objectively captured metrics to predict meeting effectiveness and inclusiveness would enable meeting measurement to scale to entire organizations. This would afford organizations the ability to understand and improve their own meetings. However, very few studies to our knowledge have attempted to do this. [93, 94] quantified the text and vocal characteristics of meeting discussions and found that certain types of conversations (e.g., conflict, social support) and expressed emotions (e.g., disappointment, excitement) are predictive of meeting success (defined in the study as a combination of factors similar to meeting effectiveness and inclusiveness). [16] analyzed body cues such as head and hand movements to predict meeting success (defined above). However, these approaches are not privacy-preserving and require participants to wear measurement technology, and therefore face challenges in scaling up to the organizational level. Moreover, they do not make use of the important meeting design characteristics that have been previously investigated. Further research

is therefore needed to explore the feasibility of predicting meeting effectiveness and inclusiveness in a privacy-preserving way using objectively measured meeting design characteristics [24].

2.5 Remote collaboration

In addition to the structural aspects of meetings, which we address in this paper, there has been significant work done on remote collaboration and the non-structural aspects of collaboration. Olson [63] provides a summary of recommendations for effective remote collaboration and best practices for remote meetings. Woolley et al. [91] studied the collective intelligence of groups and showed that groups with more equal distribution of turn-taking and groups with more equal distributions of gender had a higher collective intelligence. Lykourantzou [49] studied how personalities affect crowd-sourced teams and found that teams without a surplus of leader-type personalities exhibited less conflict and their members reported higher levels of satisfaction and acceptance. Kulkarni [42] studied massive online classes and found that the more geographically diverse the discussion groups, the better the performance of the students. Kiesler and Sproull [41] studied electronic mail systems and showed how they increased the flow of information in organizations, and in particular reduced social contexts such as location, distance, time, organizational hierarchy, age, and gender.

3 METHODOLOGY

The current work relies on subjective ratings of meeting effectiveness and inclusiveness collected via an in-client survey within a CMC system, together with meeting participation and attributes captured via telemetry (e.g., the meeting size, length, video usage, etc.). The CMC system randomly showed the survey at the end of remote meetings for the participating organizations in a pilot program. The resulting data, combined with our company's survey data, construct a valuable dataset from real meetings. We apply predictive and descriptive modeling strategies to answer the key research questions using this data.

Section 3.1 tackles RQ1 to some extent by presenting the design and implementation of a scalable survey tailored to measure the effectiveness and inclusiveness of meetings. The decisions made regarding the survey's design and deployment are crucial as they have substantial implications for data quality. We will discuss these implications in detail in Section 5, providing a detailed investigation to address all aspects of RQ1.

We provide insights into the pilot program in Section 3.2, where we maintain the confidentiality of participating organizations. Detailed information about the dataset is available in Section 3.3. In Section 3.4 we elucidate the modeling techniques employed for the analysis and extraction of insights from this dataset.

3.1 Survey Development and Implementation

In line with [24], we define meeting *effectiveness* as the extent to which business goals are attained, and meeting *inclusiveness* as the extent to which participants feel they have an opportunity to contribute and all voices have equal weight. This approach aims to minimize subjectivity by avoiding multiple perspectives or asking about decision satisfaction and applies to a broad range of workplace meetings. Moreover, it is short and simple enough to be used in large-scale deployment within organizations.

We designed a meeting effectiveness and inclusiveness survey that engaged meeting participants immediately after they left the meeting. This ensures that ratings are as proximate as possible to the actual meeting and therefore minimizes bias in recollection. However, it also means that post-meeting activities such as sharing meeting notes or actions do not impact the survey responses. Therefore, our operational definition of meeting effectiveness necessarily focuses only on activities that happened before or during the meeting.

(a) Page 1

(b) Page 2

Fig. 1. Initial End-of-Meeting Survey. We dropped the second page and only showed the first page according to randomized A/B experiment results.

Initially, our end-of-meeting survey consisted of two pages as seen in Figure 1. The first page consisted of two questions:

- (1) How effective was this meeting at achieving the business goals?
- (2) How included did you feel in the meeting?

Users could provide a rating between 1 and 5 stars, cancel out of the survey, or not provide an answer altogether, and the survey would time out after 30 seconds. The CMC system randomly selects meetings where all participants are shown the survey at the end of the call. This design provides a random sample that represents different meeting experiences in an organization. The low triggering rate of the survey reduces the chance of the same user being exposed to the survey frequently. More description about the choices of survey design is provided in section 3.1.1.

If a user gave anything less than a 5-star rating, they would be shown a second page. The user could select one or more “problem tokens” from 13 options that they felt would have improved the meeting experience. Users could provide verbatim feedback in the “Other: please specify” text box.

Like the first page, users could choose not to select anything, and the survey would time out, or users could click cancel.

We launched this two-page survey within our organization and conducted randomized experiments with the survey design, including the interface, and survey frequency. These experiments led to changes in the design, including the elimination of the second page. The final survey that is used for the current analysis contains a single page (Figure 1 (a)).

3.1.1 Experimentation: Two randomized controlled experiments were designed to compare different survey interfaces and triggering logic (frequency):

- (1) Comparing two survey interfaces: Control (A): The two-page survey as shown in Figure 1. Treatment (B): One-page survey (removed the second page of the survey shown in 1 (b)). The control survey gathers more meeting experience details but comes with a higher cognitive load, potentially reducing data quality and quantity. Comparing both surveys helps gauge the impact of survey length on data quality in this context. This experiment was run for 6 weeks. During this time, the control and treatment populations (about 8K users combined) received the respective surveys in randomly selected 10% of their meetings. The main hypothesis was that the one-page survey allows for more a diverse rating distribution (less biased and more informative data).
- (2) A two-factor experiment designed to choose “Trigger Rate”, i.e., the percentage of meetings selected for the survey, and “cool-down period”, i.e., the minimum amount of time that is needed to lapse between two consecutive survey exposures for a user. This experiment ran for 2 weeks within a population of 4.6K users. The main hypothesis was that less frequent survey exposure leads to a more diverse rating distribution (healthier and more informative data).

The four treatments were:

- (a) 1% trigger rate, no cool-down
- (b) 5% trigger rate, no cool-down
- (c) 15% trigger rate, no cool-down
- (d) 15% trigger rate, 7-day cool-down

The experiments were conducted in non-overlapping time periods. The results from the experiments led to clear conclusions that were implemented in our pilot program with external organizations:

- (1) The one-page survey enables more distributed ratings. This is demonstrated by a statistically significant drop in the responses with 5-star ratings on both questions, i.e., “Perfect Meeting Rate” (PMR).
- (2) The 7-day cool-down also corresponds to the lowest PMR.
- (3) A high trigger rate, if not accompanied by a cool-down period, can lead to more skewed ratings.

The one-page survey was finally implemented with a built-in 7-day cool-down period and a 10% trigger rate. For global organizations, additional rules were implemented to follow the local rules as needed. Importantly, the relatively low trigger rate (survey frequency) and the 7-day cool-down period both minimize the impact of the survey on user behavior during meetings, as participants are not constantly asked to complete the survey.

Telemetry collection happens entirely behind the scenes and automatically by the application without any user involvement.

Company	Employee Count (approximate)	Industry	Countries Included	Total Responses
A	32,000	Telecom	Global	2,624
B	>500,000	Consulting	US/Canada Only	2,450
C	20,000	Consulting	US/Canada Only	1,306
D	20,000	Consumer Goods	Global	971
E	5,000	Agriculture and Construction	US Only	66

Table 1. Description of participating organizations and dataset size.

3.2 Dataset and Pilot Program

After launching the survey internally, we partnered with five global companies ranging from 5,000 employees all the way to over 500K employees to build a diverse dataset. Table 1 gives an overview of the different types of organizations participating in this study. Data was collected in aggregate during the period between March 1st - June 30th, 2022, however, the exact timelines differed by participating organizations. All companies adopted hybrid policies by 2022, but the exact number of days spent in the office vs. remotely is unknown. Our data only contains ratings from users who participate in the meeting via the CMC system. Although the majority of meetings in our data are likely to be all-remote (i.e., users participate remotely via the CMC system), our data cannot identify and exclude participants that were physically co-located in hybrid meetings (see also Section 7.3 for discussion of this). After applying filters, approximately 7,330 usable responses from these organizations remained in the dataset. This accounts for 61% of the unfiltered collected ratings. Filters had to be applied to ensure that ratings were valid for our analysis. The main filters are >2 participants, call duration < 150 minutes, and time taken to complete the survey > 4 seconds. The filter on call duration is in place to exclude calls that may not confirm a similar pattern, purpose, or goal as in regular meetings. This filter has removed about 1% of all data points. The final dataset consists of both the external organization and internal pilot data: 15K ratings in total after applying filters.

3.3 Dataset Variables

Survey responses and call telemetry are linked using a shared identification number, protecting the privacy of both the respondent and the meeting participant. Call telemetry provides limited but secure insights into meeting attributes, including technical aspects and partial user behavior, with no personal identifiers. Consequently, demographic attributes are not included in this data. While this absence of personal information can be a challenge when modeling effectiveness and inclusiveness metrics, it also ensures that the resulting model and insights can be seamlessly integrated into real-world applications without requiring sensitive information.

In the rest of this section, we introduce the main variables in the current work.

3.3.1 Outcome variables. The survey used in the current work contains two questions about meeting effectiveness and inclusiveness on a 5-point star rating scale (Figure 1 (a)). We define two binary variables from these ratings: *Effective* and *Inclusive*, defined as 1 if a 4- or 5-star rating and 0 otherwise. Every other variable used in this study comes from meeting telemetry. *Effective* and *Inclusive* are the only variables in this study that are provided by users and are our main outcome variables in modeling.

We also have a telemetry-based outcome variable: Participation. Participation is 1 if the user participates vocally in the meeting and 0 otherwise. It is computed based on the Number of Encoded audio Frames (NEF) throughout the call. NEF is recorded per participant and does not require any audio recording. It is merely based on counting the number of audio packets from one participant during the call. NEF is normalized by meeting size to have the same scale as the “proportion of the meeting that the attendee spoke” (this is because audio frames are only sent when a voice activity detector is triggered). We consider $NEF > 10\%$ as a proxy for “participating in the conversation”. So Participation is 1 if $NEF > 10\%$ and 0 otherwise. Note that this does not mean exactly speaking for 10% of the meeting duration. It is a proxy for participation that is determined based on a small user study and correlation analysis with the real data. We conducted a small user study where users provided consent to share the audio content. In these calls, we compared the duration of participating in a conversation and normalized the NEF. The data showed that less than 10% of the values are the results of greetings at the beginning and end of the meeting rather than meaningful participation in the conversation. This was later confirmed by the high correlation pattern observed with other outcome metrics using the main data.

We will refer to the three variables Effective, Inclusive, and Participation as Effective Inclusive Meeting (EIM) outcome metrics.

3.3.2 Independent variables. In addition to Participation, telemetry provides detailed information on these areas: meeting duration, each attendee’s call duration, number of participants (meeting size), choice of media³ and its duration by each participant, minimal information on meeting’s scheduling metadata such as time of day, day of the week, and type⁴, whether a USB headset was in use, general statistics on network condition, and audio/video signal processing statistics.

In our initial modeling steps, we use binary variables. To transform continuous variables into binary ones, we select thresholds or ranges that show the highest sensitivity to the EIM outcome metrics when analyzed individually: We scan across a wide range of thresholds, calculate the lift⁵ in the EIM outcome variables, and choose the threshold that results in the most significant lift. In the case of ties, we opt for a middle value or one that holds practical significance

For example, we measure video usage by Video Duration Percent $> 30\%$. Since each participant in a group call can independently choose to enable or disable their video, we measure video usage on a per-participant basis. In our analysis, video duration refers to the portion of the call when a participant both viewed others’ videos and shared their own. Our initial evaluations showed that video duration is only beneficial if relative to the call duration. Hence, we define video usage as a percentage of the call duration. Video Duration Percent varies between 0% and 100%. To create a binary variable, we compute the lift in Effective and Inclusive by Video Duration Percent $> t$ for multiple values of $t \in [0.01, 0.9]$. The lift is highest when $t \approx 0.3$. Similarly, we found interesting binary variables by converting Call Duration to Short call (10 min. or less) and Meeting Size to Small Meeting (8 or less).

Meeting telemetry includes the type of the meeting: recurring (repeated regularly, such as weekly), scheduled (one-off meeting invitations), or ad-hoc (calls that people in a group chat initiate without a prior calendar invite). This study excludes ad-hoc meetings since they do not receive the Effective Inclusive Meeting (EIM) survey. We used the binary variable Recurring in modeling and insights development.

³Media refers to audio, video, or screen sharing.

⁴Scheduling frequency

⁵Lift in a binary variable X by another binary variable Y is the conditional probability of $X=True$ given $Y=True$ divided by the overall probability of $X=True$. It measures how much a change in one variable is associated with a similar change in the other variable.

We also utilized the rich telemetry about network and signal processing statistics to generate composite metrics for call quality and call reliability.

Quality Issues is a binary classifier that consumes 40 telemetry statistics about issues such as echo, noise, or speech distortions. We trained this classifier independently on ground truth from the Call Quality Feedback (CQF) survey that is displayed at the end of a random subset of calls in the CMC system⁶. Call quality ratings are the most accurate measures of call quality available from real calls. It is a single rating that reflects the user's opinion about the overall quality and is collected immediately after the call ends. *Quality Issues* is a gradient-boosting decision tree (lightGBM [40]) that is trained to predict the probability of poor CQF rating (rating 1 or 2 out of 5-stars). We measure the performance of binary classifiers by the Area Under the Curve (AUC) of the Receiver Operating Characteristic curve. This classifier has 74% AUC on the validation set. A 74% AUC is considered good performance when using call telemetry to predict real user ratings on entire call quality in this CMC system. Call telemetry is aggregated statistics and does not fully capture the changes in quality during the call. For this and other reasons, predicting user ratings just based on aggregated statistics remains a challenging task.

Reliability Issues is a simple aggregation of telemetry about call drop, one-way audio, or similar problems. This metric does not require advanced machine learning solutions since most reliability problems are well detectable by the application itself, and there is already binary telemetry to record their presence. The metric *Reliability Issues* is 1 if the call involves any reliability problem, 0 otherwise.

From a user perspective, *Reliability Issues* captures whether users can join and stay in a meeting remotely. In contrast, *Quality Issues* capture the technical audio/video *quality* of their experience if they can join and stay in the meeting.

3.4 Modeling Methodology

Our main analysis applies three modeling techniques:

- EIM graphical model [24] to detect the most important attribute of *Effective* and *Inclusive* meetings and their correlation structure
- Generalized Linear Models (GLM) [27] to explore interaction effects in sub-graphs
- Gradient-Boosting Decision Tree, such as lightGBM [40], to test the predictive power of available telemetry for EIM metrics

We apply the algorithm introduced in [24] (with minor modifications) to the three EIM outcome metrics: *Participation* → *Inclusive* → *Effective*. In this graph, each node is a variable, and each directed edge represents an adjusted Odds Ratio (OR) from a multivariate GLM model. Directed edges are from independent to outcome variables.

The algorithm to fit this model [24] has two main steps. First, the neighborhood for each node is determined using $l1$ -regularized logistic regression. This is shown to provide a close approximation of an optimum graph structure (see section 3.3.5.1 in [24]). The result of this step is a graph structure or sets of neighborhoods for each main outcome variable. The algorithm is simplified by fixing the main outcome variables to *Inclusive*, *Effective*, and *Participation*.

The hierarchy from *Participation* to *Inclusive* and then *Effective* is determined by the clues from the data and literature. We set the order between *Effective* and *Inclusive* based on the Akaike Information Criteria (AIC) [27]: model AIC value is lower when *Inclusive* predicts *Effective*. A similar result is reported in [24]. *Participation* is the third outcome node in this graph because of its significant role in connecting EIM metrics with many meeting attributes. *Participation* is not only the strongest predictor of *Inclusive*, but it also correlates with more

⁶This survey uses a Likert scale to measure the quality of the call

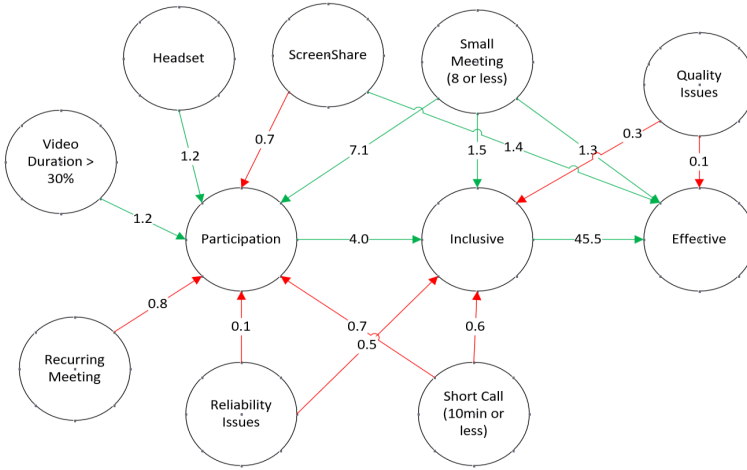


Fig. 2. The graphical model showing the network of conditional dependencies between meeting attributes and Effective and Inclusive. The red and green edges show negative and positive dependencies, respectively. The weight on each edge is the adjusted Odds Ratio as a comparable measure of strength for each dependency.

actionable attributes like video and headset usage. Without Participation as a link, we couldn't learn the detailed effects and interactions among these attributes. In addition to the support from data, this hierarchy aligns with prior research showing that participation is important for inclusiveness [4] and that inclusiveness drives meeting effectiveness [23, 60].

The second step in fitting the graph model from [24] is estimating the weights of graph edges. The coefficients of $l1$ -regularized logistic regression are not proper candidates for this purpose. While the regularization technique used in this step is an effective way of reducing the number of variables, it does not provide a valid ground for statistical hypothesis testing on these parameters [37]. Hence the second step of the algorithm applies GLM modeling without regularization to estimate the weights for each edge and further prune the graph if the parameters are statistically zero at a 95% confidence level.

In a separate modeling task in this work, we use LightGBM to fit a predictive model. LightGBM is a gradient-boosting decision tree and is less interpretable than linear models. This property makes LightGBM more appropriate and successful in predictive tasks than descriptive modeling. The set of variables that are used as input for this model is larger and includes more granular variables since there are no interpretability requirements for the predictive model.

4 EFFECTIVE INCLUSIVE MEETING (EIM) MODEL RESULTS

This section addresses RQ2 and RQ3 by constructing models based on the connections between meeting characteristics and Effective and Inclusive metrics as described in Section 3.3.1. To answer RQ2, we fit the EIM graph model using the two-step algorithm described in Section 3.4. Initially, we considered all available binary variables, allowing the algorithm to determine their relevance as predictors in the EIM graph model. Figure 2 shows the results from analyzing the entire dataset comprising 7,330 data points. Red and green edges in the figure denote negative and positive dependencies, respectively, with edge weights representing adjusted ORs computed from GLM coefficients for each node. The findings from this model and related follow-up GLM models are

discussed in Sections 4.1, 4.2, and 4.3. In addition, Section 4.4 utilizes Gradient-Boosting Decision Tree modeling to address RQ3.

4.1 Odds Ratio as Graph Weights

The **odds** of an event is the ratio between the probability of an event occurring and not occurring. For example, if the overall probability of a meeting being Effective is 90%, then the odds of it being Effective is $0.9/0.1$ or “9 to 1”.

ORs are ratios between the odds of an event under two different conditions formed by another, e.g., the rate between the odds of Effective meetings with and without Quality Issues in the call. If there is no dependency, then the OR is close to 1. OR values greater than 1 represent positive dependency and lower than 1 represent negative dependency. ORs provide a standardized measure to compare the strength or importance of attributes regardless of any assumption about the conditional rate of the outcome variable.

For example, 0.1 on the edge between Quality Issues and Effective means that the odds of an Effective meeting experience are 90% lower when the call has quality issues. Likewise, the 1.3 between Small Meeting (8 or less) and Effective means that attendees in meetings with at most 8 participants have 30% higher odds of having an Effective meeting experience. Critically, the 0.1 and 1.3 numbers are meaningful and comparable regardless of any baseline distribution for Effective and can be fairly compared with each other, hence our choice. However, if converting odds into probabilities, the results are only meaningful under an assumed rate of Effective under no Quality Issues and no Small Meeting. For example, to convert the $OR=0.1$ between Quality Issues and Effective into the % change in “the probability of Effective”, we need to specify the probability of Effective in a call without quality issues (baseline). Assuming that this baseline is 95%, $OR=0.1$ means the 95% chance of Effective drops by 66% in the presence of quality issues. Note that the 66% now is relative to the 95% baseline and would change with a different baseline. Since the baseline is specific for each attribute, the % change in probability of the model attributes is not necessarily comparable. Therefore, OR, being free of such assumptions, is the appropriate metric for interpreting multivariate models with binary variables.

4.2 Key Attributes

Inclusive and Effective have the highest correlation in the EIM graph. We believe this strong correlation is exaggerated due to common-method variance [67] and the survey response characteristics; see Section 5 for an in-depth discussion of this. The weight is large enough to encourage collapsing the two variables and creating a composite outcome variable for the model. However, our experiments showed that combining them into a single outcome variable weakens the descriptive power of the model. For example, the effects of Participation and other predictors are statistically strongest if using Inclusive as a stand-alone outcome metric. Similarly, a model that predicts Inclusive allows for a better understanding of interactions between inputs than a composite outcome variable. Therefore, we kept Inclusive and Effective as separate nodes in the graph.

After the Inclusive – Effective edge, the highest correlated pairs are:

- (1) Participation – Inclusive (OR = 4.0)
- (2) Small Meeting – Participation (OR = 7.1)
- (3) Quality Issues – Effective/Inclusive (ORs = 0.1 and 0.3)
- (4) Reliability Issues – Participation (OR = 0.1)

Section 4.3 describes a deep analysis of these areas. It is worth mentioning some of the variables that were expected to be correlated with meeting effectiveness and inclusiveness but were dropped by the algorithm. These are great examples of the importance of large data to help validate prior

theories or hypotheses about what matters in meeting inclusiveness and effectiveness. Below is a list of variables whose correlation with inclusiveness and effectiveness was not strong enough to stay in the final model.

- **Day of week:** Prior research shows that more multitasking during meetings, associated with less effective meetings, happens Mondays to Thursdays (compared to Fridays), and that Mondays are associated with high boredom levels at work [15, 52]. Our univariate analysis indeed showed slightly higher Effective and Inclusive rates on Fridays compared to Mondays. However, the effect was not strong enough to stay in the model in the presence of more dominant factors and vanished quickly in the modeling process.
- **Time of day:** Prior research shows that more multitasking during meetings happens in the mornings (compared to afternoons) and that people engage in more focused work in the afternoon [15, 52]. In our data, whether the meeting occurred in the morning or the evening proved to be irrelevant to the EIM outcome metrics. The flexible work hours during the pandemic may have caused this pattern.
- **Busy day:** Prior research shows that people multitask during meetings to catch up on their workload (including a high number of meetings) [15]. In our data, the rating distributions are not different for people who have a large number of meetings (10 or more calls) on that day vs. people with fewer meetings per day.

4.3 Insights

4.3.1 Participation: As shown in Figure 2, Participation is the most important predictor of the Inclusive and subsequently Effective nodes in this model. This is in line with qualitative and survey research showing that meeting participation is key for meeting satisfaction and overall employee engagement (reviewed in [4]). Participating in conversation is associated with 4x higher odds of having an Inclusive experience. This is equivalent to an 8% higher probability of having an Inclusive experience for attendees whose rate of having an Inclusive experience is 35% to 60% when not participating. Interestingly, this delta is smaller for attendees with baselines below 35% or above 60%. Figure 3 shows the details of this change. The figure displays a lower impact of participation when the baseline for inclusiveness is not already very high or very low. When the Inclusive rate is already quite high, it is not surprising that participation has limited scope for impact. However, when the Inclusive rate is very low, this suggests that factors other than Participation may be more important (e.g., the tone of meetings or other aspects of team culture [4]). It is crucial for organizations to know their baseline Inclusive rate before starting any campaign for improving meeting culture.

Insight: Vocally participating in meetings is associated with the largest change in the probability of having an inclusive experience (8% increase) for meetings with a mid-range baseline probability of being Inclusive (35-60%).

The strong link between Participation and Inclusive (and thereby Effective) motivated us to analyze the extent to which *all* attendees participated in a meeting (rather than an individual) and how this correlates with attendees' experiences (see Section 4.3.6).

4.3.2 Meeting Size: The second strongest correlation belongs to Meeting Size. The model predicts that attendees in meetings with less than eight people have respectively 50% and 30% higher odds of rating meetings as Inclusive and Effective. Also, they are more likely to participate by a large margin (7 times higher odds).

Our findings extend earlier qualitative and survey research on meeting effectiveness [5, 20, 34]. Large meetings may reflect situations in which some invited participants are not relevant to the meeting, leading to real or perceived inefficiencies in time use among some or all participants

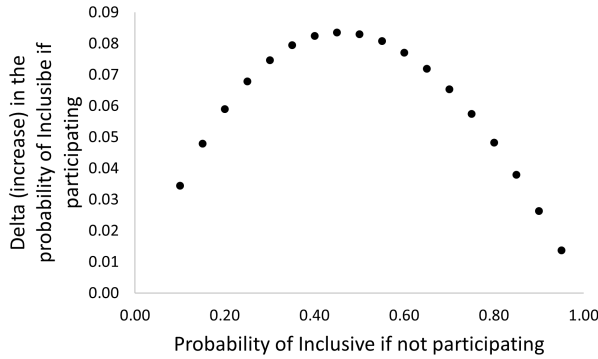


Fig. 3. Model predictions for the impact of Participation in the probability of Inclusive experience for different Inclusive rate baselines. The impact declines for participants who would have had a very high Inclusive rate (over 60%, already very good experience) or a very low Inclusive rate (lower than 35%, significantly poor experience) without participating in conversations.

[20, 34, 75]. Larger meetings also provide fewer opportunities to participate (see Figure 4). This is in harmony with the findings in [75]. Given the aforementioned link between participation and inclusiveness, it is not surprising that larger meetings are also associated with lower perceived inclusiveness.

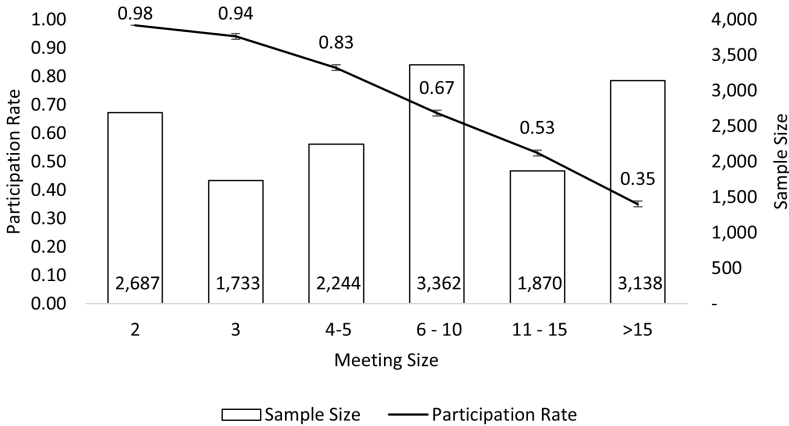


Fig. 4. Participation rate drops as Meeting Size increases. The rapid decline begins with more than 5 people in the meeting. The difference in Participation is largest and statistically significant when we compare 8-or-less-participant meetings with more-than-8-participant meetings.

We sought to better understand the link between Meeting Size and Effective and Inclusive meetings under different scenarios. Specifically, we were interested in the linear effect of meeting size (i.e., rather than a binary variable defining small meetings as 8 attendees or less) to gain a more granular understanding. We were also interested in how this varied by whether a meeting was part of a recurring series or whether it was a one-off meeting, as these are broadly associated with different meeting purposes (e.g., team stand-ups vs. brainstorming meetings) that may be

differently impacted by meeting size [4]. Additionally, we looked at the role of meeting duration, as prior research shows that long meetings are perceived as less effective [4]. To do these analyses, we fit separate GLM models to predict Effective and Inclusive meetings using Meeting Size in its numeric format (not binary), with Recurring and Call Duration as inputs. Given that the above aspects can all interact in meaningful ways (e.g., a long and large recurring meeting may have a lower effectiveness than a long and large one-off meeting), we include interaction terms in the model. GLM parameters for the Effective outcome variable are available in Table 7 in the appendix. The results demonstrate lower Effective and Inclusive rates for larger meetings. GLM predictions under different scenarios show that if we increase the Meeting Size from 2 to 14, then the Inclusive rate falls from about 98% to 94% and, the Effective rate falls from 97% to 93%. However, this decline is not identical for all types of meetings. The model estimates the largest negative delta for recurring meetings that take no longer than 30 minutes; for these meetings, every two new participants are associated with a reduction in the chance of an Inclusive or Effective experience by 1% absolute.

Insight: Meetings with fewer attendees are associated with much higher odds of vocal participation and of being rated as Effective and Inclusive. The impact of meeting size on effectiveness and inclusiveness is strongest for short recurring meetings (up to 30 minutes).

One prominent example of a short, recurring meeting is the daily stand-up meeting. Small-scale studies of daily stand-ups in software teams suggest that Meeting Size is indeed an important factor in the perceived success of such meetings [84, 85]. Here we corroborate and generalize these findings to a range of teams and industries, and identify specific quantitative effects of meeting size that are actionable.

4.3.3 Call Quality and Reliability: The third group of dominant factors in the model is the quality and reliability of the call (the Quality Issues and Reliability Issues nodes). Reliability issues occur twice more often than quality issues in this data, with strong connections to the Participation and Inclusive nodes. Given that reliability refers to basic task completion (merely being able to participate in conversation), we expect it to be necessary for a good meeting experience. But is it sufficient? The EIM graphical model demonstrates an interesting order between Quality Issues and Reliability Issues. To better reveal this pattern, we show the relevant graph weights in Table 2, with rows and columns indicating connecting nodes. It shows that as we move from the basic task completion (Participation) to ultimately Effective and Inclusive meetings, we move from a high Reliability correlation to a balanced mix with Quality and finally to just Quality as a highly relevant factor.

Insight: Whereas call reliability is critical for enabling participation and an inclusive meeting experience, call quality becomes essential to achieve both inclusive and effective meetings.

	Participation	Inclusive	Effective
Quality Issues	-	0.35	0.14
Reliability Issues	0.13	0.49	-

Table 2. Pattern of Quality and Reliability weights for connecting nodes in the graphical model. Rows and columns indicate connecting nodes. The empty cells (-) show that there is no edge between the corresponding two nodes in the graph.

4.3.4 Media: A remote participant has different ways of participating in meetings (i.e., different media choices):

- (1) Audio-only: only speaking via the microphone

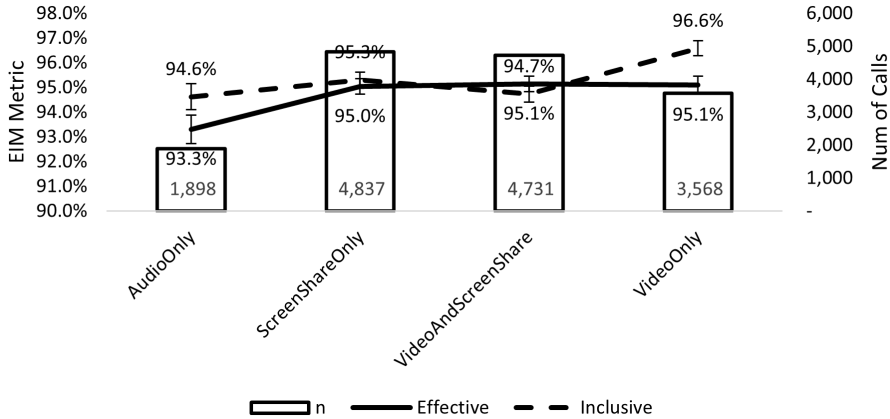


Fig. 5. EIM metrics by media. The least Effective and Inclusive experiences are audio-only calls. Video-only calls are most Inclusive (2% higher than audio-only), but not the most Effective. 95% Confidence intervals are shown for each metric to enable comparisons.

- (2) Video: turning on their webcam video or seeing other attendees' videos
- (3) Screen-sharing: sharing their screen or seeing other attendees' screens

About 87% of calls in this data involve participants using video or screen-sharing for various portions of the call duration. How does meeting inclusiveness and effectiveness vary with different media choices? Figure 5 shows effective and inclusive ratings for different media choices. We find that audio-only calls are the least Inclusive and the least Effective, while the most Inclusive calls are video calls, where such calls are 2% (absolute) more Inclusive than audio-only calls.

Comparing video and audio-only calls shows that the Participation rate grows when video duration increases. To analyze this further, we look at the proportion of the call duration that involved video, defining Video Duration Percent as the rate of video duration divided by call duration. For example, if a participant joins a meeting for 30 minutes and uses video for 15 minutes, we see a 50% Video Duration Percent. We quantify the effect of Video Duration Percent on Participation by the lift in the probability of Participation. We find that for an average meeting, as video duration changes from less than 10% of the call to greater than 70%, the lift in the probability of participation increases from 8% to 11%. In other words, longer video duration is associated with a larger increase in the probability of participation.

While this data demonstrates video as one of the attributes of inclusive meetings, it also indicates that not all calls may benefit from this feature. The lift in Participation depends on the Meeting Size and Call Duration. To investigate this, we used a GLM model that predicts Participation probability by Meeting Size and Call Duration. We determined the thresholds for meeting size (8 attendees) and video duration (30 minutes) by repeating a similar sensitivity analysis that we used to generate binary variables out of continuous ones for the graph models. Table 8 in the Appendix contains the details of the GLM model. We find that using video is associated with the largest lift in Participation rate in short meetings with few participants: in a meeting with less than eight participants that takes 30 minutes or less, the chance of participation can increase by more than 6% if attendees use video for at least 30% of call duration. For larger meetings, the associated change in participation can even become negative when the meeting duration passes 40 minutes.

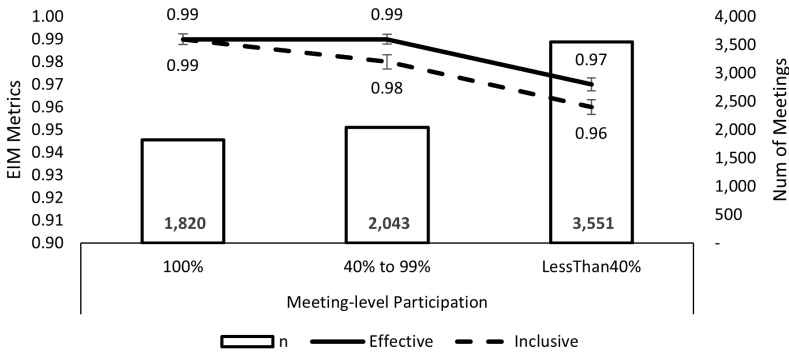


Fig. 6. EIM metrics by Participation rate per meeting. Meeting with less than 40% participation are up to 2% less Inclusive and Effective

ScreenShare is the only attribute with positive and negative correlations in the EIM graph (with Effective and Participation nodes). ScreenShare is a feature that allows a participant to present her screen during the meeting or see other people's screens. Meetings with screen-sharing have 30% lower odds of Participation while 40% higher odds of Effective rating. Meeting Size is partly explaining this pattern: meetings that have at least 10% of their time on screen-sharing are 50% larger than meetings with less or no screen-sharing on average. Additionally, during meetings with screen-sharing, the presenter generally talks more than the audience, resulting in lower overall participation across meeting participants.

Insight: Longer video duration by meeting participants is associated with a larger increase in the probability of participation. Using video in meetings is most strongly associated with an increase in the probability of participation in short meetings (30 minutes or less) with few participants (less than eight participants). Meetings with screen-sharing are associated with higher odds of being rated as effective, and with a lower rate of participation (at least partly due to these meetings tending to be larger in size, and involving presentations in which presenters speak more than other participants).

Our results support the notion that video plays a significant role in establishing a social presence, the “sense of being with another” [11], which is associated with multiple positive communication outcomes such as trust and enjoyment [62]. Other studies also report that video is primarily effective in smaller groups where it can afford a sense of intimacy among participants [8, 9].

4.3.5 Headset: According to the graph model, the odds of Participation increase by 20% if the user is using a headset. Further analysis shows that a specific group of calls drives this correlation: audio-only or screen share-only calls, i.e., calls without video. Audio-only calls with headsets have about 5% higher Participation rates than audio-only calls without a headset.

It is worth mentioning that detecting headset usage from telemetry is a challenging task because many devices report incorrect telemetry. Therefore, we expect our “no-headset” label to have contaminated some headset users. For this reason, the current predictions may underestimate the effect of the headset on participation in no-video meetings.

4.3.6 Per-Meeting Participation Rate. All previous analyses on Participation and its effect on the rate of Inclusive meetings are about correlations between attendees' ratings and their call telemetry. It shows that if one person participates in the conversations, there is a higher chance that they rate meetings as Inclusive. However, previous research suggests that the choice to

participate in meeting conversations may be influenced by social norms set by other participants [9]. We investigated this possibility and looked at the impact of *every attendee* participating on effectiveness and inclusiveness ratings. This requires aggregating both Participation and survey ratings for each meeting.

Participation data is available for all meeting attendees, but ratings are usually provided by at most one of the meeting attendees due to the participants' non-response⁷: only 3% of the meetings have more than one rating. Therefore, in our dataset, meeting-level ratings are not very different from user-level ratings. However, whenever there is more than one rating available for a meeting, we use their average as the meeting-level rating of effectiveness and inclusiveness. On the other hand, we use three bins to aggregate Participation at the meeting level. These are the three bins for per-meeting Participation metric:

- 100% (all attendees participated): a proxy for “everyone spoke more than once”
- 40% - 99% of attendees participated
- Less than 40% of attendees participated

Figure 6 shows how both Effective and Inclusive rates decline together with the Participation rate across the three bins. We find that meetings where everyone participates are 1.8% more Effective and 2.6% more Inclusive than meetings where everyone doesn't participate. By contrast, meetings with less than 50% participation are 1.5% less Effective and 2.7% less Inclusive than meetings with more or equal to 50% participation.

Insight: The extent of participation *across participants within a meeting* is associated with meeting effectiveness and inclusiveness. Meetings where everyone participates are rated as more effective and inclusive than those where this isn't the case.

This confirms that the choice of Participation for an attendee not only impacts her own meeting experience but also can influence the collective experience of all attendees in the meeting. This aligns with research suggesting that the choice to participate is influenced by social norms [9], where observing participation from others may create a safe space for oneself to participate, thereby increasing the perception of inclusiveness, and ultimately effectiveness. Indeed, Constantinides et al. [23] found that perceiving one's peers to be “comfortable sharing their thoughts and making contributions” is related to one's comfort to contribute in a meeting; both are aspects of inclusiveness or psychological safety.

4.4 Generalizability and Predictive Power

In this section, beyond identifying a set of meeting attributes that are associated with meeting effectiveness and inclusiveness, we address two key questions that each determine the wider value of our approach and findings. Firstly, we test whether the EIM graph model is generalizable across organizations with different teams, practices, and aims, and therefore different meeting experiences. This is crucial because it tests the wider generalizability not only of our meeting-related insights but also of our measurement and modeling approach more broadly. Secondly, we test whether the available telemetry data we consider is sufficient to reliably *predict* the survey responses to the effectiveness and inclusiveness questions for specific attendees. This is important because such predictability would enable avoiding surveys altogether, thereby minimizing any potential observer effect (see also 7.3) and organizational burden, and enabling a wider deployment of our approach at scale.

⁷Non-response refers to missing data in surveys where respondents decide not to fill out or submit their answers. In the current research, everyone in the meeting receives the survey, so all attendees with no ratings are considered a non-response.

Target	Input	Combined	Subset 1	Subset 2
Effective	Quality Issues	0.14 (<0.01)	0.1 (<0.01)	0.2 (<0.01)
Effective	Inclusive	45.48 (<0.01)	54 (<0.01)	39.7 (<0.01)
Effective	ScreenShare	1.39 (<0.01)	1.3 (0.05)	1.39 (0.01)
Effective	Small Meeting (8 or less)	1.29 (<0.01)	-	1.43 (<0.01)
Inclusive	Quality Issues	0.35 (<0.01)	-	0.25 (<0.01)
Inclusive	Reliability Issues	0.49 (<0.01)	0.44 (<0.01)	0.47 (<0.01)
Inclusive	Participation	4.05 (<0.01)	4.47 (<0.01)	3.58 (<0.01)
Inclusive	Small Meeting (8 or less)	1.51 (<0.01)	1.43 (0.02)	1.56 (<0.01)
Inclusive	Short Call (10min. or less)	0.61 (<0.01)	0.5 (<0.01)	-
Participation	Reliability Issues	0.13 (<0.01)	0.1 (<0.01)	0.16 (<0.01)
Participation	Recurring	0.82 (<0.01)	0.84 (<0.01)	0.8 (<0.01)
Participation	ScreenShare	0.71 (<0.01)	0.72 (<0.01)	0.71 (<0.01)
Participation	Small Meeting (8 Or Less)	7.13 (<0.01)	7.39 (<0.01)	7.03 (<0.01)
Participation	Short Call (10min. or less)	0.72 (<0.01)	0.66 (<0.01)	0.79 (0.03)
Participation	Headset	1.16 (<0.01)	1.16 (0.01)	1.16 (0.01)
Participation	Video Duration > 30%	1.17 (<0.01)	1.17 (0.01)	-

Table 3. Graph parameters are adjusted ORs computed from GLM coefficients. We show the p-values in parentheses next to each parameter. The empty cells (-) show that there is no edge between the corresponding two nodes in that graph. Subset 1 and Subset 2 are mutually exclusive subsets of data in a way that they include non-overlapping participating organizations while maintaining a similar sample size.

4.4.1 EIM Graph Generalizability. The purpose of the EIM graph model is its descriptive capacity to provide a data-driven structure of correlated attributes of Effective and Inclusive meeting experiences. The insights from this model are reliable and applicable only if this structure does not alter dramatically for a new organization (i.e., the extent of model generalizability). To examine this, we split the data into two different subsets (Subset 1 and Subset 2), defined by different participating organizations, such that the two subsets have a comparable sample size and therefore comparable statistical power. We then fit the EIM graph to these two different subsets and compared the results.

Table 3 shows the three model parameters side-by-side. The parameters of our main graph (Figure 2) are under the “Combined” column. The parameters of graphs fitted to the two subsets are listed in columns “Subset 1” and “Subset 2”. These graphs have a few edges less than the Combined graph. This makes the Combined model a meta-graph of these two subsets. This is the result of more edges being pruned (not passing the statistical significance test) during the graph modeling on the two subsets. Other than that, the algorithm suggests exactly the same structure in all three cases: There is no neighboring node of Participation, for example, that would move to the Effective or Inclusive neighborhood. Therefore there is no sign indicating that the graph model and insights derived from it cannot be generalized.

The absence of some edges (correlations) in the two subsets is most likely due to the lack of statistical power in sub-samples. This is also demonstrated by the increase in p-values for common parameters when we move from the subset graphs to the Combined model. This pattern emphasizes the importance of sample size and the danger of relying on a small sample for detecting the drivers of meeting experiences.

4.4.2 Predictive Power for EIM Metrics. Our second question concerns the power of the available telemetry data to predict the survey ratings of meeting effectiveness and inclusiveness. Our ultimate aim with predictive modeling is to obviate the need for surveys altogether, so as to reduce potential

	AUC +/- error
Effective	0.65+/-0.02
Inclusive	0.72+/-0.02

Table 4. Cross-validation AUC on 50 randomly selected test sets for models predicting Inclusive than Effective ratings using available telemetry only.

observer effects and organizational burdens and enable wider deployment of our measurement and modeling approach. Predictive modeling has a different goal and requirements than descriptive modeling. Because our aim is to obviate the need for surveys, our predictive model should not include any survey-based feature data. Moreover, not being bound by descriptive purposes, there is more flexibility about the type of the model and its complexity when choosing a predictive model (e.g., no constraints on structure). We, therefore, trained two separate models to predict Effective and Inclusive ratings without Inclusive as an input variable.

In this study, the prediction task is a binary classification where a lightGBM model predicts the probability of an Effective (or Inclusive) rating. All available telemetry and some engineered features from telemetry are predictors (i.e., independent variables). These add up to 40 different inputs for each model. Table 9 provides the details about these features.

To measure the predictive power of each model, we use the Area Under the Curve (AUC) of the Receiver Operating Characteristic curve. We follow a random sampling cross-validation strategy where we repeatedly split the data into random train and test subsets, and record the model AUC on the test set. Table 4 contains the cross-validation result of 50 random samplings. The results show that there is more power in predicting Inclusive than Effective ratings. However, neither of the models demonstrates a high enough AUC for reliably predicting individual ratings in the absence of a subjective survey. We should note that the unit of random splits is single ratings and not any unique identifier of users. Therefore, it is possible that ratings from the same person on two different meetings fall in train and test, which is a potential data leak that can cause an overestimation of AUC, i.e., the current AUC is the best value that this data can provide. We can still safely conclude that meeting telemetry by itself is not adequate to predict individual ratings.

Next, we measure the change in AUC when the model moves to predict ratings for meetings outside the organizations in the training set. This is especially critical for being able to use the model in a widely deployed CMC system to automatically predict ratings for meetings without survey ratings. To test the model performance in this scenario, we run the cross-validation while we keep an entire organization outside the training and test data. We repeat this process with three different organizations. The results show a small but statistically significant drop in AUC for both Effective and Inclusive models (see Table 5). That means the models' accuracy in predicting survey ratings will reduce slightly if the model is used to predict ratings in an organization that is not part of the training set. Note that the average AUC on training cross-validation in Table 5 is lower than in Table 4 because the training dataset in the former analysis of generalizability is smaller and less diverse.

Predicting individual effectiveness and inclusiveness ratings is statistically more challenging than deriving common patterns of correlations. The former requires more data points and a broader set of predictive variables. Our study shows that meeting telemetry can provide a robust understanding of the common factors related to effectiveness and inclusiveness. However, without further research, telemetry alone cannot replace a subjective survey by predicting individual ratings. This emphasizes the importance and, by extension, the quality of the survey-based measurement

of meeting effectiveness and inclusiveness; Section 5 below details our work to understand and improve the survey data.

Unseen Test	Cross-validation AUC when training		Unseen test set AUC	
	Inclusive	Effective	Inclusive	Effective
Test organization 1	0.68+/-0.01	0.60/-0.01	0.66	0.59
Test organization 2	0.70+/-0.01	0.63/-0.01	0.66	0.58
Test organization 3	0.69+/-0.01	0.62/-0.01	0.67	0.60
Average	0.69	0.62	0.67	0.59

Table 5. AUC when moving the model to an unseen test set.

5 SURVEY MEASUREMENT OF MEETING EFFECTIVENESS AND INCLUSIVENESS: SURVEY SKEW AND ITS SOLUTIONS

The survey-based measurement of meeting effectiveness and inclusiveness is core to our methodology as, to our knowledge, it provides the only tractable way of measuring these subjective and complex constructs (e.g., [23]). The survey data, therefore, serves as the “ground truth” in our descriptive and predictive modeling, with the value of telemetry attributes depending heavily on the quality of the survey data. Data quality becomes even more important as surveys are deployed in the real-world organizational context. Ensuring high-quality survey data is therefore a key focus. One of the main challenges we observed is participants’ tendency to provide 4 or 5-star ratings regardless of their true experience (i.e., what is termed “skew” in the data). Table 6 shows the relative frequencies of star ratings provided by survey respondents; a significant portion of surveys is rated with 4- or 5-star ratings. This section describes studies we conducted to understand the skew in ratings and approaches to mitigate its impact.

This skew pattern has been reported in other settings that rely on user ratings, such as online marketplaces like eBay, Uber, and AirBnB [33, 86]. There are several potential complementary explanations for this. First, low ratings could incur a perceived reputation cost for yourself and others: low ratings could encourage retaliatory behavior from those being rated or could incur other social costs, particularly in a workplace context [86]. Raters may therefore avoid low ratings out of fear of this cost to themselves or others. Sections 5.1 and 5.6 describe analyses addressing this. Second, due to the ubiquity of star rating systems and aforementioned reputation costs, participants may have learned norms that responding with 4 or 5 stars is the appropriate response regardless of context (see Section 5.1 for more on this) [33]. Third, and related to the previous reasons, participants may have limited capacity to respond to the survey, so defaulting to the norm of 4 or 5 stars is the fastest and safest approach in terms of reputation and cognitive costs of providing informative ratings (see Sections 5.2, 5.3, and 5.4 for more on this).

5.1 Replacing stars with worded labels

One hypothesis is that perceived reputation costs associated with low ratings discourage such ratings (e.g., raters don’t want to disparage the meeting organizer), leading participants to provide 4- or 5-star ratings [86]. This may be further exacerbated by the ubiquity of online star rating systems which reinforce learned norms around how to respond [33]. [33] showed that replacing a star rating system with positive-skewed worded options (response options that have more positive- than negative-valenced verbal labels) substantially reduces rating skew in an online marketplace. First, a positive-skewed option set may reduce the perceived reputation cost (for yourself and others)

	Effective	Inclusive
1-star	1%	1%
2-star	1%	1%
3-star	3%	3%
4-star	12%	10%
5-star	82%	86%

Table 6. Relative frequency of star ratings provided by respondents shows a high tendency to rate 4 or 5 stars for meeting Effectiveness and Inclusiveness.

of providing lower ratings; second, worded options may discourage respondents from relying on learned norms around star ratings. We conducted a small experiment to test this approach in the meeting rating context. A random subset of 2,000 employees was polled via email and randomly allocated to receive a link to one of two surveys as displayed in Figure 7 (we relied on email as a rapid approach to initial experimentation). The *control* survey invited the respondents to rate the effectiveness and inclusiveness of their last meeting on a scale of 1 to 5 stars, and the *treatment* survey instead asked respondents to rate the meeting on a scale of “Not Effective”, “Somewhat Effective”, “Effective”, “Quite Effective”, “Very Effective” (with analogous options for inclusiveness).

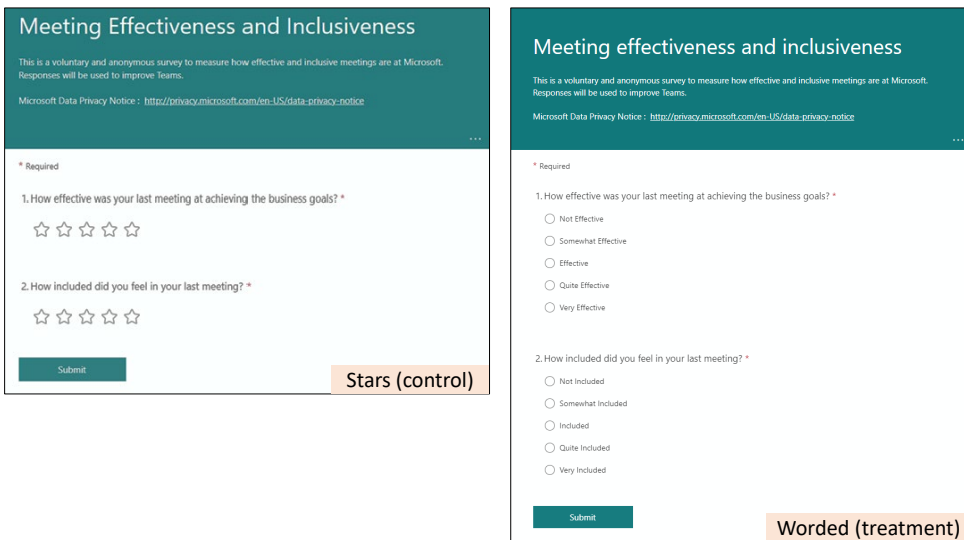


Fig. 7. Participants received either a control survey where they rated meeting effectiveness and inclusiveness out of 5 stars or a treatment survey where they rated instead out of positively skewed worded options.

We had a response rate of 18.1% and 17.2% for the control and treatment conditions, respectively (no significant difference between these rates using a Fisher’s exact test, $p = 0.639$). When comparing the distribution of responses between conditions, we found a substantial difference for both effectiveness and inclusiveness ratings, as seen in Figure 8. A Fisher’s exact test of the counts across ratings showed a statistically significant difference between conditions for both rating questions ($p < 0.001$ for both). To further estimate the improvement with the worded options, we computed the Shannon Entropy in each condition as a measure of the information in the data. We observed

an increase from 1.44 to 2.07 bits for effectiveness (a 44% increase), and an increase from 1.48 to 2.04 for inclusiveness (a 38% increase), respectively. Thus, we demonstrate as a proof-of-concept that positive-skewed worded response options can decrease rating skew, relative to star ratings, for ratings of meeting effectiveness and inclusiveness. This is a promising direction and should be further tested in the CMC system when feasible.

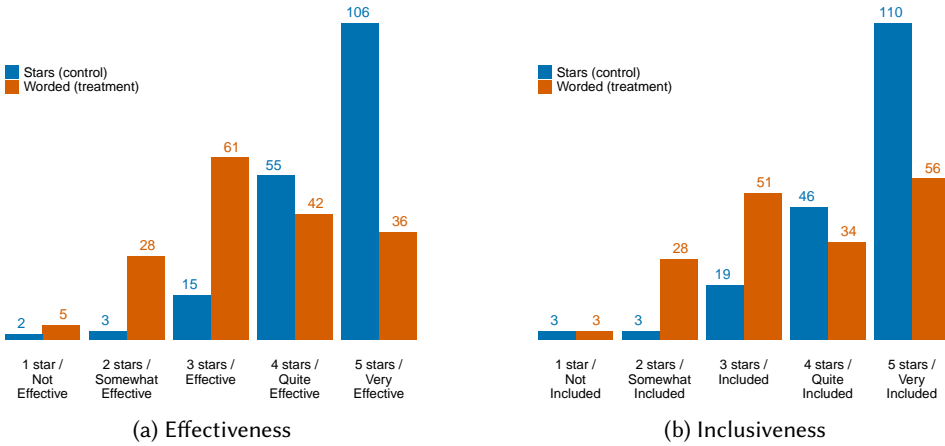


Fig. 8. Rating skew experiment results for ratings of meeting (a) effectiveness and (b) inclusiveness. Figures show the distribution of responses for the stars (control) condition and the worded (treatment) condition. Numbers indicate the counts of participants selecting that option.

5.2 Survey completion time

Another hypothesis for the cause of the skew was that respondents were responding with 4 or 5 stars too quickly to get through the survey (i.e., they had limited capacity to respond) and those responses weren't reflective of their experience. The work done in [90] examined the appropriateness of using response speed as an indicator of data quality. We used their recommended filter of 4 seconds for the 2-question survey. Comparing the skew on the data with the filter of response time, the results showed a reduction of users rating 5 on both inclusive and effectiveness by 14 percentage points absolute. The rate of 4 or 5 on inclusiveness and effectiveness individually were both reduced by 3%. Since this filter was able to improve data quality by reducing the skew, it was applied to the data before the analysis and modeling presented in Section 4.

5.3 Survey fatigue

As the survey was available for a longer period of time we noticed an increase in the rate of 4- and 5-star responses, as well as a decrease in the survey response rate. One suspected reason for these changes was user fatigue with the survey. A longitudinal study was conducted by grouping the data based on how long it had been since the user had previously experienced the survey. Then we could look at the rate of responding 5 to both questions for the different groups. There is a bias where users who have more meetings are more likely to observe the survey and have shorter times between survey exposures. To address this bias, we did this analysis comparing people with similar survey exposure counts. What we found was a critical point where 7 days after receiving the survey there was a decrease in users rating 5 to both questions by 1.6%. This supports having a

cool-down period between survey requests. This was implemented in our CMC system for this work and executed as detailed in Section 3.1.1.

5.4 Timing of meetings

Another scenario potentially contributing to the skew was the timing of the meeting for the user. We identified for these studies the time of day the meeting occurred if the user had a meeting after the meeting occurred, and if they had multiple meetings that day.

We did a study on the subset of the data where the time from when the current meeting ended to when the next meeting began was within 2 hours. On this data, we compared the response rate and metric rates for different times until the next meeting began. We found that when the time to the next meeting is less than 5 minutes from when the questionnaire is shown there is a 2% lower response rate, and the users were more likely to rate the meeting effectiveness and inclusiveness lower. This directly shows this is not a cause of the rating skew. Moreover, there is no need to implement a filter that suppresses the survey in the event of users having meetings close to each other in the CMC system.

Considering ratings by the time of day that the meeting occurred, we found that there was no statistical difference in the rating distribution.

The last thing we checked in terms of someone having a busy day was the number of meetings they had using the CMC system that day. What we found was the more meetings someone had in a day the less likely they were to respond to the survey, and the more likely they were to rate the meeting 5 stars. We looked for an interaction effect for the time of day for people with large volumes of meetings, but the difference in metric distribution for this group was not significant.

5.5 User demographics

We hypothesized there is a demographic landscape of raters that we are not covering in our modeling. People with different jobs are expected to have different rating behaviors. We defined user cohorts based on net user behavior, looking at how the individual engaged with the CMC system over the past month. To do this, we used the total number of meetings, the average meeting size, the frequency at which the user hosts meetings, and the percentage of meetings rated. There are multiple ways to cluster data. Our priority was to have a clustering that could give us information about different demographics. To maintain information about the clusters we partitioned each feature into 2 bins based on a set of thresholds. These thresholds are 56 for the number of meetings, 30% for the percentage of meetings rated, 20% for the percentage of hosted meetings, and 10 for the average meeting size. Looking at all possible combinations of these and the resulting metrics there was a clear trend to have 3 cohorts defined as:

- Cohort 0: The percentage of rated meetings is low.
- Cohort 1: The percentage of rated meetings is high and the number of meetings is low.
- Cohort 2: The percentage of rated meetings is high and the meeting size is large.

These cohorts had statistically different distributions of the results. With Cohort 0 showing the least skew, and Cohort 2 showing the most. These cohorts proved to be useful in modeling. However, we were not able to use them for this study given that the cohort computation requires historical data and many of the organizations in our data were relatively new. Future work should explore this further.

5.6 User anonymity

User perceptions of meeting inclusiveness and effectiveness are sensitive data to collect. Users might be wary of answering due to the potential impact on themselves or their coworkers. To test

if anonymity was important we conducted an experiment with surveys that were sent via email instead of integrated into the CMC system. Users were sent emails with a link to a survey, we had two surveys with the same questions but with different privacy statements. The first statement mirrored the survey through the CMC system, providing a link to a data protection notice. The second statement had a stronger worded privacy statement that additionally included that their responses were anonymous. We found no statistical difference in the rating distribution of the two surveys. This analysis did not provide evidence that letting raters know their responses were anonymous would resolve the skew.

6 SUMMARY OF FINDINGS

In this work, we developed descriptive and predictive models for meeting effectiveness and inclusiveness and showed that they generalize across multiple organizations. We conducted an analysis using the models and found several novel insights about the attributes driving meeting effectiveness and inclusiveness (see Insights in Section 4.3). We also analyzed and shared solutions for the survey skew which is the major data quality challenge in the subjective measurement of meeting effectiveness and inclusiveness.

We summarize the results of our research in three areas related to the three research questions described in Section 1.

- CMC systems can implement automated randomized surveys about meeting effectiveness and inclusiveness to show at the end of a meeting. These surveys, deployed at a low triggering rate, can produce a steady flow of subjective measurements. The resulting data is reliable for measuring the overall patterns and changes in meeting experiences over time. The main challenge is data quality, particularly the survey rating skew. We can address rating skew by using worded positively-skewed response options and by implementing a delay time between two consecutive survey exposures for a user. Additionally, considering a minimum response time per question is a powerful approach for post hoc filtering of potentially invalid survey response data. For this purpose, the CMC system should submit the survey completion time alongside other telemetry.
- The descriptive model and data confirm that participating in conversations is the most significant factor in meeting inclusiveness and effectiveness, both at the individual and group levels. Beyond that, small meetings are also associated with increased participation. Meetings with ScreenShare are more effective but not necessarily more inclusive. On the other hand, video usage can help with inclusiveness for small and short meetings.
- Statistical models are beneficial for decomposing the subjective ratings of effectiveness and inclusiveness onto telemetry-based attributes. However, their predictive power is limited in predicting individual ratings based on the meeting telemetry data available today alone. The descriptive models are robust and reliable concerning the shifts in the underlying data (e.g., across organizations), but the moderate predictive power of statistical models can suffer significantly when applied to unseen data.

7 DISCUSSION

7.1 Theoretical implications

Understanding and improving meetings has been a central focus of meeting science since the pioneering work of Schwartzman [79]. The vast majority of such research has relied on small-scale surveys that often ask about people's general meeting experiences, rather than enabling analysis at the meeting level (e.g., [20, 43]), or involves small-scale analyses of meetings using manual coding of behavior (e.g., [38, 45]). The COVID-19-related shift to increased remote and hybrid work

has provided an opportunity to understand and improve meetings at scale using CMC systems. Here we demonstrate the feasibility and value of a large-scale, cross-organizational approach to measuring and linking subjective survey data on meeting effectiveness and inclusiveness, together with objective telemetry data on meeting participation and other attributes captured via a CMC system during real-world remote meetings.

We posit that large-scale measurement of meetings is key to an in-depth understanding of the factors that contribute to successful meetings. Moreover, it enables organizations to contextually understand their own meetings and rigorously test their own policy-based or technological interventions for improving them, thereby addressing the heterogeneity common to behavioral change [14]. Indeed, this value is recognized by the organizational executives we partnered with, who wanted a clear line of sight to our model's ability to show statistically significant changes for organizational interventions to improve meeting culture.

Objective data on meeting participation and other attributes (i.e., telemetry) is central to our approach as it can ultimately enable at-scale and passive measurement of meeting experiences [22]. Meeting effectiveness, and especially inclusiveness, are both complex constructs with many subjective elements that are difficult to currently capture using objective measures (e.g., [57, 71]); hence, we relied on survey ratings to measure these and conducted in-depth investigations to improve survey data quality. However, the ultimate aim of our approach for real-world deployment is to build a predictive model using objective telemetry that can accurately estimate meeting effectiveness and inclusiveness without the need for survey ratings. To this end, we demonstrate that our descriptive model generalizes well across organizations, but that predicting effectiveness and inclusiveness in the absence of survey ratings requires further work to increase the range of telemetry included in the model, and to improve the quality of survey data (see also Section 7.3).

7.2 Design implications

7.2.1 Improving meeting design and technologies. Our analyses of meeting attributes point to opportunities for improving meeting design and technologies. For example, given the importance of participation for meeting effectiveness and inclusiveness, there is a need to reduce barriers to participation. To this end, recent work has developed a detector for failed speech interruption attempts which can help participants take the floor in conversations [31]. Similarly, meeting dashboards may also provide participants with insights about their conversations to encourage more equitable participation [77]. To address the negative impact of meeting size, there is an opportunity to nudge meeting organizers during scheduling to help them reflect on their intended participants, their workflow, and meeting goals, and decide whether meeting size can be reduced [55].

Our findings underscore the importance of video use for participation and ultimately inclusiveness in small meetings. Ensuring that all participants can transmit reliable and high-quality video is a top priority for improving CMC systems. Moreover, given the fact that video use is strongly influenced by social norms—people's decision to turn on video is at least partly determined by whether others in the meeting also do so [8–10]—there is an opportunity to use in-system reminders to encourage video use in relevant meeting contexts. Alternatively, as a mitigation against video fatigue [10], there is an opportunity to explore whether avatars can serve a similar purpose in increasing social presence and meeting participation [64, 92] (see also Section 7.3).

7.2.2 Deployment of large-scale meeting measurement systems. Deploying large-scale measurement and modeling of behavior into organizations invites understandable concerns and requires buy-in. We encountered several challenges in the setup of the study. First, survey respondents needed more clarity about how responses would be used, reflecting common concerns about workplace

surveillance (which can include hierarchical as well as peer-to-peer surveillance [6, 82]), particularly in the context of the automated collection of telemetry data. If deployed inappropriately and without employee buy-in, workplace surveillance has the potential to reduce beliefs in organizational fairness, trust in leadership, and commitment [17], as well as increase stress [69]. Considering the (explicit and perceived) purpose, invasiveness, frequency and regularity, and transparency of large-scale meeting measurement systems are essential prior to wider deployment [68]. Secondly, as mentioned above, organizational executives wanted confidence in the models' ability to show statistically significant changes associated with changes they made to meeting culture. Relatedly, there were concerns about the model's explainability to external stakeholders. Indeed, explainability is important for justifying decisions (particularly if meeting measurement systems become tied to employee performance metrics), and increasing understanding and therefore control of how a system operates [3]. Lastly, during survey piloting, external organizational stakeholders in Human Resources and IT were not comfortable with sending company-wide surveys via mass emails. Alternatives like posting the survey to internal social channels resulted in a minimal response rate (less than 3%), preventing our ability to collect a diverse baseline using this methodology. Sufficient lead time to acquire organizational buy-in, wider considerations of employee rights and perceptions, and iterative testing and feedback are therefore important for the successful deployment of such systems.

7.3 Limitations and future research

Although we have explored multiple important meeting attributes, there remain further opportunities to improve the model's predictive power by including richer telemetry data. First, participation in remote meetings also includes the use of chat and reactions, which are used in ways that can contribute to effectiveness (e.g., by providing easy ways to share relevant information without disrupting ongoing conversations) and inclusiveness (e.g., by widening participation opportunities for those that may not be able or feel comfortable to participate vocally) [78].

Second, prior survey research has found that the impact of participation on meeting outcomes depends on the content and context of participation: meeting citizenship behaviors, such as sharing helpful information or ideas, can improve perceived effectiveness and overall engagement, whereas counterproductive behaviors, such as criticizing others or complaining, can harm it [2, 46, 61]. Those with meeting-relevant knowledge are more likely to participate and thereby perceive their meetings as being more effective, particularly in meeting contexts with high participant disagreement [48]. For privacy reasons, our data does not have information on the content of verbal participation, yet recent research shows that such content may be particularly predictive of meeting effectiveness and inclusiveness [93, 94]. Future work should therefore consider how to leverage such data in a privacy-preserving manner (see also below).

Third, the dynamics of participation, including patterns in turn-taking and other aspects of conversation flow, can yield rich insights into effectiveness and inclusiveness in a privacy-preserving manner (i.e., without considering the content of speech) [31, 51, 70]. Potentially fruitful data here includes patterns in the timing and duration of speech across participants [51], choral responses like laughter [13], acoustic features like prosody [57], gestures such as head nods [77], and eye gaze patterns [26].

Fourth, as mentioned above, avatars and other mixed reality technologies provide alternative media choices for participants to convey their presence in meetings without relying on video [64]. Understanding the impact of these choices on effectiveness and inclusiveness, and the interaction with video use will be important.

Relatedly, hybrid meetings (with both on-site and remote attendees) pose another challenge and opportunity for our methodology which has so far focused on remote attendance (though the

majority of meetings in our data are likely to be all-remote, our data cannot identify participants that were co-located in hybrid meetings). The user experience for post-meeting surveys requires careful consideration about where and when to deliver surveys for on-site meeting attendees, ensuring that survey responding is easy and privacy-preserving for everyone. Additionally, telemetry requires further processing to accommodate multiple on-site attendees (e.g., speaker diarization to accurately measure vocal participation for each person). Given the social and interaction asymmetries common to hybrid meetings, measuring meeting effectiveness and inclusiveness for such meetings is a priority [76].

Fifth, there is an opportunity to integrate telemetry about the presence of agendas or other pre-meeting materials, action points, and other post-meeting minutes that are known to be key for meeting effectiveness [4, 23, 55]. For example, a lack of clear goals, agendas, and post-meeting summaries have been shown to be negatively correlated with meeting quality [20] and meeting effectiveness [32, 34, 43].

Lastly, we expect that meeting type (whether a meeting is for brainstorming, decision-making, etc.) has a significant role in modeling meeting effectiveness and inclusiveness [4]. However, telemetry does not contain such information, and current machine learning (ML) solutions cannot reliably extract it automatically from the meeting invite. Making this information available has a high potential for improving EIM model accuracy.

While telemetry provides opportunities for more accurate measurements, it also imposes limitations on the scope of available attributes. There are aspects (such as participant demographics) that are expected to have statistical predictive power but are not available to a CMC system due to privacy and security concerns. This can prevent the model's ability to achieve 100% descriptive or predictive power of individual subjective assessments of effectiveness and inclusiveness. However, it also enables any ML solution based on this data to be more secure and fair in real-world applications.

As explored in Section 5, one major challenge in the subjective measurement of meeting effectiveness and inclusiveness is reducing the commonly observed skew of ratings. Survey design, such as the design of response options and the timing of deployment, is key for mitigating this; data can be refined further by applying appropriate filters (e.g., based on response times). Other opportunities include using nudges based on social norms or incentives to encourage participants to respond more frequently and honestly [18, 81]. Future work should test these ideas in new contexts and at scale, and explore other ways of improving data quality.

The "observer effect" poses an additional potential challenge to our methodology [50]. Participants' meeting interactions and survey responses may be influenced by their awareness of being observed and measured via surveys and telemetry. Although we cannot exclude this possibility, several lines of evidence argue against it. Firstly, as discussed in Section 5.6, informing participants that their survey responses would be anonymous, thereby reducing the expectation of them being specifically observed, did not influence their ratings. Secondly, other than during the study introduction and the post-meeting surveys, participants had no indication that meetings were being measured (including via telemetry, which includes only standard data regularly logged for engineering purposes). Thirdly, we chose a relatively low survey frequency (10% of meetings) and a 7-day cool-down period to minimize the impact of the survey on user behavior. Moreover, neither participants nor their managers were provided with the results of the survey ratings or telemetry during the study period, thereby precluding the influence of feedback on performance and subsequent ratings. Lastly, our observed skewed distribution of survey ratings is similar to that observed in online platforms where users likely do not have strong expectations of being observed by experimenters (at least during normal use) [86]. Similarly, given that the study was conducted over four months in a real-world context, participants may have habituated to the study,

thereby reducing any potential observer effects. Future work could further estimate the extent of the observer effect in the current context, for example, by testing whether pre-meeting or pre-study information about subsequent measurement influences participant behavior.

Lastly, our methodology requires participants to voluntarily complete the surveys, which opens up the potential for a self-selection bias in the study sample. That is, the people who took the time to complete the surveys may be systematically different from the overall population of employees (e.g., they may be less busy and/or more engaged in the workplace than the overall population). Indeed, as Section 5.5 suggests, there are differences in rating behavior based on factors such as how many meetings a person organizes and the size of their meetings. The privacy-preserving nature of our methodology precluded us from capturing more detailed demographics. However, future work could estimate and mitigate self-selection bias by capturing relevant demographic variables about participants, comparing them to the demographics of the overall population, and targeting surveys accordingly to minimize bias. Ultimately, improving the modeling of meeting effectiveness and inclusiveness using telemetry will decrease the need for self-selected survey data.

ACKNOWLEDGMENTS

We thank Scott Inglis, Thierry Tremblay, and Dejan Ivkovic for their valuable work that enabled survey development and data collection for this research.

REFERENCES

- [1] Joseph A. Allen, Tammy Beck, Cliff W. Scott, and Steven G. Rogelberg. 2014. Understanding workplace meetings: A qualitative taxonomy of meeting purposes. *Management Research Review* 37, 9 (Aug. 2014), 791–814. <https://doi.org/10.1108/MRR-03-2013-0067>
- [2] Joseph A. Allen, Nale Lehmann-Willenbrock, and Nicole Landowski. 2014. Linking pre-meeting communication to meeting effectiveness. *Journal of Managerial Psychology* 29, 8 (Nov. 2014), 1064–1081. <https://doi.org/10.1108/JMP-09-2012-0265>
- [3] Amina Adadi and Mohammed Berrada. 2018. Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI). *IEEE Access* 6 (2018), 52138–52160. <https://doi.org/10.1109/ACCESS.2018.2870052>
- [4] Joseph A. Allen and Nale Lehmann-Willenbrock. 2022. The key features of workplace meetings: Conceptualizing the why, how, and what of meetings at work. *Organizational Psychology Review* (Sept. 2022), 204138662211292. <https://doi.org/10.1177/20413866221129231>
- [5] Joseph A. Allen, Jiajin Tong, and Nicole Landowski. 2020. Meeting effectiveness and task performance: meeting size matters. *Journal of Management Development* ahead-of-print, ahead-of-print (March 2020). <https://doi.org/10.1108/JMD-12-2019-0510>
- [6] Mark Andrejevic. 2002. The Work of Watching One Another: Lateral Surveillance, Risk, and Governance. *Surveillance & Society* 2, 4 (Sept. 2002). <https://doi.org/10.24908/ss.v2i4.3359>
- [7] Tanachia Ashikali and Sandra Groeneveld. 2013. Diversity Management in Public Organizations and Its Effect on Employees' Affective Commitment. *Review of Public Personnel Administration* (Nov. 2013). <https://doi.org/10.1177/0734371X13511088>
- [8] Karolina Balogova and Duncan Brumby. 2022. How Do You Zoom?: A Survey Study of How Users Configure Video-Conference Tools for Online Meetings. In *2022 Symposium on Human-Computer Interaction for Work*. ACM, Durham NH USA, 1–7. <https://doi.org/10.1145/3533406.3533408>
- [9] Nancy Baym, Rachel Bergmann, Adam Coleman, Ricardo Reyna Fernandez, Sean Rintel, Abigail Sellen, and Tiffany Smith. 2021. Collaboration and Meetings. In *The New Future of Work: Research from Microsoft on the Impact of the Pandemic on Work Practices*. Microsoft. <https://www.microsoft.com/en-us/research/publication/collaboration-and-meetings/>
- [10] Andrew A. Bennett, Emily D. Campion, Kathleen R. Keeler, and Sheila K. Keener. 2021. Videoconference fatigue? Exploring changes in fatigue after videoconference meetings during COVID-19. *Journal of Applied Psychology* 106, 3 (March 2021), 330–344. <https://doi.org/10.1037/apl0000906>
- [11] Frank Biocca, Chad Harms, and Judee K. Burgoon. 2003. Toward a More Robust Theory and Measure of Social Presence: Review and Suggested Criteria. *Presence: Teleoperators and Virtual Environments* 12, 5 (Oct. 2003), 456–480. <https://doi.org/10.1162/105474603322761270>

- [12] Nicholas Bloom. 2022. Hybrid is the Future of Work. *[Stanford Institute for Economic Policy Research (SIEPR): Stanford, CA, USA (2022)]*.
- [13] F. Bonin, N. Campbell, and C. Vogel. 2012. *CogInfoCom 2012: 3rd IEEE International Conference on Cognitive Infocommunications*. IEEE, Piscataway, N.J. OCLC: 835886329.
- [14] Christopher J. Bryan, Elizabeth Tipton, and David S. Yeager. 2021. Behavioural science is unlikely to change the world without a heterogeneity revolution. *Nature Human Behaviour* 5, 8 (July 2021), 980–989. <https://doi.org/10.1038/s41562-021-01143-3>
- [15] Hancheng Cao, Chia-Jung Lee, Shamsi Iqbal, Mary Czerwinski, Priscilla N Y Wong, Sean Rintel, Brent Hecht, Jaime Teevan, and Longqi Yang. 2021. Large Scale Analysis of Multitasking Behavior During Remote Meetings. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. ACM, Yokohama Japan, 1–13. <https://doi.org/10.1145/3411764.3445243>
- [16] Jun-Ho Choi, Marios Constantinides, Sagar Joglekar, and Daniele Quercia. 2021. KAIROS: Talking Heads and Moving Bodies for Successful Meetings. In *Proceedings of the 22nd International Workshop on Mobile Computing Systems and Applications*. ACM, Virtual United Kingdom, 30–36. <https://doi.org/10.1145/3446382.3448361>
- [17] Rebecca M. Chory, Lori E. Vela, and Theodore A. Avtgis. 2016. Organizational Surveillance of Computer-Mediated Workplace Communication: Employee Privacy Concerns and Responses. *Employee Responsibilities and Rights Journal* 28, 1 (March 2016), 23–43. <https://doi.org/10.1007/s10672-015-9267-4>
- [18] Adrienne Chung and Rajiv N. Rimal. 2016. Social norms: A review. *Review of Communication Research* 4 (2016), 1–28. <https://doi.org/10.12840/issn.2255-4165.2016.04.01.008>
- [19] Scott Cliff, Joseph A. Allen, Steven G. Rogelberg, and Alex Kello. 2015. Five theoretical lenses for conceptualizing the role of meetings in organizational life. In *The Cambridge handbook of meeting science*. New York, NY: Cambridge University Press, 20–46.
- [20] Melissa A. Cohen, Steven G. Rogelberg, Joseph A. Allen, and Alexandra Luong. 2011. Meeting design characteristics and attendee perceptions of staff/team meeting quality. *Group Dynamics: Theory, Research, and Practice* 15, 1 (2011), 90–104. <https://doi.org/10.1037/a0021549>
- [21] Marios Constantinides, Sagar Joglekar, and Daniele Quercia. 2021. Retrofitting Meetings for Psychological Safety. <http://arxiv.org/abs/2109.12976> arXiv:2109.12976 [cs].
- [22] Marios Constantinides and Daniele Quercia. 2022. The Future of Hybrid Meetings. In *2022 Symposium on Human-Computer Interaction for Work*. ACM, Durham NH USA, 1–6. <https://doi.org/10.1145/3533406.3533415>
- [23] Marios Constantinides, Sanja Šćepanović, Daniele Quercia, Hongwei Li, Ugo Sassi, and Michael Eggleston. 2020. ComFeel: Productivity is a Matter of the Senses Too. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 4, 4 (Dec. 2020), 1–21. <https://doi.org/10.1145/3432234>
- [24] Ross Cutler, Yasaman Hosseinkashi, Jamie Pool, Senja Filipi, Robert Aichner, Yuan Tu, and Johannes Gehrke. 2021. Meeting Effectiveness and Inclusiveness in Remote Collaboration. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW1 (April 2021), 1–29. <https://doi.org/10.1145/3449247>
- [25] Owen Daly-Jones, Andrew Monk, and Leon Watts. 1998. Some advantages of video conferencing over high-quality audio conferencing: fluency and awareness of attentional focus. *International Journal of Human-Computer Studies* 49, 1 (July 1998), 21–58. <https://doi.org/10.1006/ijhc.1998.0195>
- [26] Ziedune Degutyte and Arlene Astell. 2021. The Role of Eye Gaze in Regulating Turn Taking in Conversations: A Systematized Review of Methods and Findings. *Frontiers in Psychology* 12 (April 2021), 616471. <https://doi.org/10.3389/fpsyg.2021.616471>
- [27] Annette J. Dobson. 2002. *An introduction to generalized linear models* (2nd ed ed.). Chapman & Hall/CRC, Boca Raton.
- [28] Amy Edmondson. 1999. Psychological Safety and Learning Behavior in Work Teams. *Administrative Science Quarterly* 44, 2 (June 1999), 350–383. <https://doi.org/10.2307/2666999>
- [29] Paul Van Eecke and Raquel Fernández. 2016. On the Influence of Gender on Interruptions in Multiparty Dialogue. 2070–2074. <https://doi.org/10.21437/Interspeech.2016-951>
- [30] Bernardo M. Ferdman and Barbara Deane (Eds.). 2014. *Diversity at work: the practice of inclusion*. Jossey-Bass, A Wiley Brand, San Francisco, CA.
- [31] Szu-Wei Fu, Yaran Fan, Yasaman Hosseinkashi, Jayant Gupchup, and Ross Cutler. 2022. Improving Meeting Inclusiveness using Speech Interruption Analysis. In *Proceedings of the 30th ACM International Conference on Multimedia*. ACM, Lisboa Portugal, 887–895. <https://doi.org/10.1145/3503161.3548379>
- [32] Boni García, Micael Gallego, Francisco Gortázar, and Antonia Bertolino. 2019. Understanding and estimating quality of experience in WebRTC applications. *Computing* 101, 11 (Nov. 2019), 1585–1607. <https://doi.org/10.1007/s00607-018-0669-7>
- [33] Nikhil Garg and Ramesh Johari. 2019. Designing Informative Rating Systems: Evidence from an Online Labor Market. <http://arxiv.org/abs/1810.13028> arXiv:1810.13028 [cs].

- [34] Jennifer L. Geimer, Desmond J. Leach, Justin A. DeSimone, Steven G. Rogelberg, and Peter B. Warr. 2015. Meetings at work: Perceived effectiveness and recommended improvements. *Journal of Business Research* 68, 9 (Sept. 2015), 2015–2026. <https://doi.org/10.1016/j.jbusres.2015.02.015>
- [35] Zixiu Guo, John D’Ambra, Tim Turner, Huiying Zhang, and Tong Zhang. 2006. Effectiveness of Meeting Outcomes in Virtual vs. Face-to-Face Teams: A Comparison Study in China. In *Americas Conference on Information Systems*. 13.
- [36] Tony F. Habash. 1999. The impact of audio- or video-conferencing and group decision tools on group perception and satisfaction in distributed meetings. *The Psychologist-Manager Journal* 3, 2 (1999), 211–230. <https://doi.org/10.1037/h0095872>
- [37] T Hastie, R Tibshirani, and J Friedman. 2001. *The Elements of Statistical Learning*. Springer, New York, NY, USA.
- [38] Marcella Hoogeboom and Celeste Wilderom. 2015. Effective Leader Behaviors in Regularly Held Staff Meetings: Surveyed vs. Videotaped and Video-Coded Observations. In *The Cambridge Handbook of Meeting Science* (1 ed.), Joseph A. Allen, Nale Lehmann-Willenbrock, and Steven G. Rogelberg (Eds.). Cambridge University Press, 381–412. <https://doi.org/10.1017/CBO9781107589735.017>
- [39] Deborah James and Janice Drakich. 1993. Understanding gender differences in amount of talk: A critical review of research. In *Gender and conversational interaction*. Oxford University Press, New York, NY, US, 281–312.
- [40] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. 2017. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems* 30 (2017).
- [41] Sara Kiesler and Lee Sproull. 1992. Group decision making and communication technology. *Organizational Behavior and Human Decision Processes* 52, 1 (June 1992), 96–123. [https://doi.org/10.1016/0749-5978\(92\)90047-B](https://doi.org/10.1016/0749-5978(92)90047-B)
- [42] Chinmay Kulkarni, Julia Cambre, Yasmine Kotturi, Michael S. Bernstein, and Scott R. Klemmer. 2015. Talkabout: Making Distance Matter with Small Groups in Massive Classes. In *Proceedings of the 18th ACM Conference on Computer Supported Cooperative Work & Social Computing*. ACM, Vancouver BC Canada, 1116–1128. <https://doi.org/10.1145/2675133.2675166>
- [43] Desmond J. Leach, Steven G. Rogelberg, Peter B. Warr, and Jennifer L. Burnfield. 2009. Perceived Meeting Effectiveness: The Role of Design Characteristics. *Journal of Business and Psychology* 24, 1 (March 2009), 65–76. <https://doi.org/10.1007/s10869-009-9092-6>
- [44] Campbell Leaper and Melanie M. Ayres. 2007. A Meta-Analytic Review of Gender Variations in Adults’ Language Use: Talkativeness, Affiliative Speech, and Assertive Speech. *Personality and Social Psychology Review* 11, 4 (Nov. 2007), 328–363. <https://doi.org/10.1177/1088868307302221>
- [45] Nale Lehmann-Willenbrock and Joseph A. Allen. 2014. How fun are your meetings? Investigating the relationship between humor patterns in team interactions and team performance. *Journal of Applied Psychology* 99, 6 (Nov. 2014), 1278–1287. <https://doi.org/10.1037/a0038083>
- [46] Nale Lehmann-Willenbrock, Joseph A. Allen, and Dain Belyeu. 2016. Our love/hate relationship with meetings: Relating good and bad meeting behaviors to meeting outcomes, engagement, and exhaustion. *Management Research Review* 39, 10 (Oct. 2016), 1293–1312. <https://doi.org/10.1108/MRR-08-2015-0195>
- [47] Nale Lehmann-Willenbrock, Stephenson J. Beck, and Simone Kauffeld. 2016. Emergent Team Roles in Organizational Meetings: Identifying Communication Patterns via Cluster Analysis. *Communication Studies* 67, 1 (Jan. 2016), 37–57. <https://doi.org/10.1080/10510974.2015.1074087>
- [48] Isaac A. Lindquist, Emily E. Adams, and Joseph A. Allen. 2020. If I Had Something to Add, I Would: Meeting Topic Competences and Participation. *Journal of Personnel Psychology* 19, 2 (April 2020), 86–96. <https://doi.org/10.1027/1866-5888/a000255>
- [49] Ioanna Lykourantzou, Angeliki Antoniou, Yannick Naudet, and Steven P. Dow. 2016. Personality Matters: Balancing for Personality Types Leads to Better Outcomes for Crowd Teams. In *Proceedings of the 19th ACM Conference on Computer-Supported Cooperative Work & Social Computing (CSCW ’16)*. Association for Computing Machinery, New York, NY, USA, 260–273. <https://doi.org/10.1145/2818048.2819979>
- [50] Ritch Macefield. 2007. Usability Studies and the Hawthorne Effect. *Journal of Usability Studies* 2, 3 (May 2007), 145 – 154.
- [51] Eleni Margariti, Sean Rintel, Brendan Murphy, and Abigail Sellen. 2022. Automated mapping of competitive and collaborative overlapping talk in video meetings.. In *CHI Conference on Human Factors in Computing Systems Extended Abstracts*. 8.
- [52] Gloria Mark, Shamsi T. Iqbal, Mary Czerwinski, and Paul Johns. 2014. Bored Mondays and focused afternoons: the rhythm of attention and online activity in the workplace. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, Toronto Ontario Canada, 3025–3034. <https://doi.org/10.1145/2556288.2557204>
- [53] Roberto J. Mejias. 2007. The Interaction of Process Losses, Process Gains, and Meeting Satisfaction Within Technology-Supported Environments. *Small Group Research* 38, 1 (Feb. 2007), 156–194. <https://doi.org/10.1177/1046496406297037>
- [54] Microsoft. 2022. *2022 Work Trend Index: Annual Report*. Technical Report. Microsoft. <https://www.microsoft.com/en-us/worklab/work-trend-index/hybrid-work-is-just-work>.

- [55] Joseph E. Mroz, Joseph A. Allen, Dana C. Verhoeven, and Marissa L. Shuffler. 2018. Do We Really Need Another Meeting? The Science of Workplace Meetings. *Current Directions in Psychological Science* 27, 6 (Dec. 2018), 484–491. <https://doi.org/10.1177/0963721418776307>
- [56] Nicol L. Davidson. 2013. *Trust and Member Inclusion as Communication Factors to Foster Collaboration in Globally Distributed Teams*. Ph.D. Dissertation.
- [57] Oliver Niebuhr, Ronald Böck, and Joseph A. Allen. 2021. On the Sound of Successful Meetings: How Speech Prosody Predicts Meeting Performance. In *Companion Publication of the 2021 International Conference on Multimodal Interaction*. ACM, Montreal QC Canada, 240–248. <https://doi.org/10.1145/3461615.3485412>
- [58] Karin Niemantsverdriet and Thomas Erickson. 2017. Recurring Meetings: An Experiential Account of Repeating Meetings in a Large Organization. *Proceedings of the ACM on Human-Computer Interaction* 1, CSCW (Dec. 2017), 1–17. <https://doi.org/10.1145/3134719>
- [59] Carol T. Nixon and Glenn E. Littlepage. 1992. Impact of meeting procedures on meeting effectiveness. *Journal of Business and Psychology* 6, 3 (1992), 361–369. <https://doi.org/10.1007/BF01126771>
- [60] J. F. Nunamaker, Alan R. Dennis, Joseph S. Valacich, Douglas Vogel, and Joey F. George. 1991. Electronic meeting systems. *Commun. ACM* 34, 7 (July 1991), 40–61. <https://doi.org/10.1145/105783.105793>
- [61] Isabelle Odermatt, Cornelius J. König, Martin Kleinmann, Maria Bachmann, Heiko Röder, and Patricia Schmitz. 2018. Incivility in Meetings: Predictors and Outcomes. *Journal of Business and Psychology* 33, 2 (April 2018), 263–282. <https://doi.org/10.1007/s10869-017-9490-0>
- [62] Catherine S. Oh, Jeremy N. Bailenson, and Gregory F. Welch. 2018. A Systematic Review of Social Presence: Definition, Antecedents, and Implications. *Frontiers in Robotics and AI* 5 (Oct. 2018), 114. <https://doi.org/10.3389/frobt.2018.00114>
- [63] Judith S. Olson and Gary M. Olson. 2013. Working Together Apart: Collaboration over the Internet. *Synthesis Lectures on Human-Centered Informatics* 6, 5 (Nov. 2013), 1–151. <https://doi.org/10.2200/S00542ED1V01Y201310HCI020> Publisher: Morgan & Claypool Publishers.
- [64] Payod Panda, Molly Jane Nicholas, Mar Gonzalez-Franco, Kori Inkpen, Eyal Ofek, Ross Cutler, Ken Hinckley, and Jaron Lanier. 2022. AllTogether: Effect of Avatars in Mixed-Modality Conferencing Environments. In *2022 Symposium on Human-Computer Interaction for Work*. ACM, Durham NH USA, 1–10. <https://doi.org/10.1145/3533406.3539658>
- [65] Jone L. Pearce and Amy E. Randel. 2004. Expectations of organizational mobility, workplace social inclusion, and employee job performance. *Journal of Organizational Behavior* 25, 1 (Feb. 2004), 81–98. <https://doi.org/10.1002/job.232>
- [66] Angie Pendergrass. 2019. Inclusive scientific meetings: Where to start. <https://opensky.ucar.edu/islandora/object/manuscripts%3A983/datastream/PDF/view>
- [67] Philip M. Podsakoff, Scott B. MacKenzie, Jeong-Yeon Lee, and Nathan P. Podsakoff. 2003. Common method biases in behavioral research: A critical review of the literature and recommended remedies. *Journal of Applied Psychology* 88, 5 (2003), 879–903. <https://doi.org/10.1037/0021-9010.88.5.879>
- [68] Daniel M. Ravid, David L. Tomczak, Jerod C. White, and Tara S. Behrend. 2020. EPM 20/20: A Review, Framework, and Research Agenda for Electronic Performance Monitoring. *Journal of Management* 46, 1 (Jan. 2020), 100–126. <https://doi.org/10.1177/0149206319869435>
- [69] Daniel M Ravid, Jerod C White, Dave L Tomczak, Ahleah F Miles, and Tara S Behrend. 2020. A Meta-Analysis of the Effects of Digital Surveillance of Workers: A Psychology Focused Approach. *Microsoft New Future of Work Symposium* (2020).
- [70] Andrew Reece, Gus Cooney, Peter Bull, Christine Chung, Bryn Dawson, Casey Fitzpatrick, Tamara Glazer, Dean Knox, Alex Liebscher, and Sebastian Marin. 2022. Advancing an Interdisciplinary Science of Conversation: Insights from a Large Multimodal Corpus of Human Speech. <http://arxiv.org/abs/2203.00674> arXiv:2203.00674 [cs].
- [71] Bruce A Reinig. 2003. Toward an Understanding of Satisfaction with the Process and Outcomes of Teamwork. *Journal of Management Information Systems* 19, 4 (April 2003), 65–83. <https://doi.org/10.1080/07421222.2003.11045750>
- [72] Darryl B. Rice, Nicole C. J. Young, and Sharon Sheridan. 2020. Improving employee emotional and behavioral investments through the trickle-down effect of organizational inclusiveness and the role of moral supervisors. *Journal of Business and Psychology* (Jan. 2020). <https://doi.org/10.1007/s10869-019-09675-2>
- [73] E. Sean Rintel. 2010. Conversational management of network trouble perturbations in personal videoconferencing. In *Proceedings of the 22nd Conference of the Computer-Human Interaction Special Interest Group of Australia on Computer-Human Interaction*. ACM, Brisbane Australia, 304–311. <https://doi.org/10.1145/1952222.1952288>
- [74] Steven Rogelberg, Desmond Leach, Peter Warr, and Jennifer Burnfield. 2006. "Not Another Meeting!" Are Meeting Time Demands Related to Employee Well-Being? *The Journal of applied psychology* 91 (Feb. 2006), 83–96. <https://doi.org/10.1037/0021-9010.91.1.83>
- [75] N.C. Romano and J.F. Nunamaker. 2001. Meeting analysis: findings from research and practice. In *Proceedings of the 34th Annual Hawaii International Conference on System Sciences*. IEEE, Maui, HI, USA, 13 pp. <https://doi.org/10.1109/HICSS.2001.926253>

- [76] Banu Saatçi, Roman Rädle, Sean Rintel, Kenton O'Hara, and Clemens Nylandstedt Klokmose. 2019. Hybrid Meetings in the Modern Workplace: Stories of Success and Failure. In *Collaboration Technologies and Social Computing*, Hideyuki Nakanishi, Hironori Egi, Irene-Angelica Chounta, Hideyuki Takada, Satoshi Ichimura, and Ulrich Hoppe (Eds.). Vol. 11677. Springer International Publishing, Cham, 45–61. https://doi.org/10.1007/978-3-030-28011-6_4 Series Title: Lecture Notes in Computer Science.
- [77] Samiha Samrose, Daniel McDuff, Robert Sim, Jina Suh, Kael Rowan, Javier Hernandez, Sean Rintel, Kevin Moynihan, and Mary Czerwinski. 2021. MeetingCoach: An Intelligent Dashboard for Supporting Effective & Inclusive Meetings. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. ACM, Yokohama Japan, 1–13. <https://doi.org/10.1145/3411764.3445615>
- [78] Advait Sarkar, Sean Rintel, Damian Borowiec, Rachel Bergmann, Sharon Gillett, Danielle Bragg, Nancy Baym, and Abigail Sellen. 2021. The promise and peril of parallel chat in video meetings for work. In *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*. ACM, Yokohama Japan, 1–8. <https://doi.org/10.1145/3411763.3451793>
- [79] Helen B. Schwartzman. 1989. *The Meeting*. Springer US, Boston, MA. <https://doi.org/10.1007/978-1-4899-0885-8>
- [80] Abigail Sellen. 1995. Remote Conversations: The Effects of Mediating Talk With Technology. *Human-Computer Interaction* 10, 4 (Dec. 1995), 401–444. https://doi.org/10.1207/s15327051hci004_2
- [81] Michael G. Smith, Maryam Witte, Sarah Rocha, and Mathias Basner. 2019. Effectiveness of incentives and follow-up on increasing survey response rates and participation in field studies. *BMC Medical Research Methodology* 19, 1 (Dec. 2019), 230. <https://doi.org/10.1186/s12874-019-0868-8>
- [82] Lebene Richmond Soga, Bernd Vogel, Ana Margarida Graça, and Kofi Osei-Frimpong. 2021. Web 2.0-enabled team relationships: an actor-network perspective. *European Journal of Work and Organizational Psychology* 30, 5 (Sept. 2021), 639–652. <https://doi.org/10.1080/1359432X.2020.1847183>
- [83] Willem Standaert, Steve Muylle, and Amit Basu. 2016. An empirical study of the effectiveness of telepresence as a business meeting mode. *Information Technology and Management* 17, 4 (Dec. 2016), 323–339. <https://doi.org/10.1007/s10799-015-0221-9>
- [84] Viktoria Stray, Dag I.K. Sjøberg, and Tore Dybå. 2016. The daily stand-up meeting: A grounded theory study. *Journal of Systems and Software* 114 (April 2016), 101–124. <https://doi.org/10.1016/j.jss.2016.01.004>
- [85] Viktoria Gulliksen Stray, Yngve Lindsjorn, and Dag I.K. Sjøberg. 2013. Obstacles to Efficient Daily Meetings in Agile Development Projects: A Case Study. In *2013 ACM / IEEE International Symposium on Empirical Software Engineering and Measurement*. IEEE, Baltimore, Maryland, 95–102. <https://doi.org/10.1109/ESEM.2013.30>
- [86] Steven Tadelis. 2016. Reputation and Feedback Systems in Online Platform Markets. *Annual Review of Economics* 8, 1 (Oct. 2016), 321–340. <https://doi.org/10.1146/annurev-economics-080315-015325>
- [87] John C. Tang and Ellen Isaacs. 1992. Why do users like video?: Studies of multimedia-supported collaboration. *Computer Supported Cooperative Work (CSCW)* 1, 3 (Sept. 1992), 163–196. <https://doi.org/10.1007/BF00752437>
- [88] Maria del Carmen Triana, Bradley L. Kirkman, and Maria Fernanda Wagstaff. 2012. Does the Order of Face-to-Face and Computer-Mediated Communication Matter in Diverse Project Teams? An Investigation of Communication Order Effects on Minority Inclusion and Participation. *Journal of Business and Psychology* 27, 1 (March 2012), 57–70. <https://doi.org/10.1007/s10869-011-9232-7>
- [89] Elizabeth S. Veinott, Judith Olson, Gary M. Olson, and Xiaolan Fu. 1999. Video helps remote work: speakers who need to negotiate common ground benefit from seeing each other. In *Proceedings of the SIGCHI conference on Human factors in computing systems the CHI is the limit - CHI '99*. ACM Press, Pittsburgh, Pennsylvania, United States, 302–309. <https://doi.org/10.1145/302979.303067>
- [90] Dustin Wood, P. D. Harms, Graham H. Lowman, and Justin A. DeSimone. 2017. Response Speed and Response Consistency as Mutually Validating Indicators of Data Quality in Online Samples. *Social Psychological and Personality Science* 8, 4 (2017), 454–464. <https://journals.sagepub.com/doi/abs/10.1177/1948550617703168>
- [91] Anita Williams Woolley, Christopher F. Chabris, Alex Pentland, Nada Hashmi, and Thomas W. Malone. 2010. Evidence for a Collective Intelligence Factor in the Performance of Human Groups. *Science* 330, 6004 (Oct. 2010), 686–688. <https://doi.org/10.1126/science.1193147>
- [92] Mika Yasuoka, Marko Zivko, Hiroshi Ishiguro, Yuichiro Yoshikawa, and Kazuki Sakai. 2022. Effects of Digital Avatar on Perceived Social Presence and Co-presence in Business Meetings Between the Managers and Their Co-workers. In *Collaboration Technologies and Social Computing*, Lung-Hsiang Wong, Yugo Hayashi, Cesar A. Collazos, Claudio Alvarez, Gustavo Zurita, and Nelson Baloian (Eds.). Vol. 13632. Springer International Publishing, Cham, 83–97. https://doi.org/10.1007/978-3-031-20218-6_6 Series Title: Lecture Notes in Computer Science.
- [93] Ke Zhou, Marios Constantinides, Luca Maria Aiello, Sagar Joglekar, and Daniele Quercia. 2021. The Role of Different Types of Conversations for Meeting Success. *IEEE Pervasive Computing* 20, 4 (Oct. 2021), 35–42. <https://doi.org/10.1109/MPRV.2021.3115879>

- [94] Ke Zhou, Marios Constantinides, Sagar Joglekar, and Daniele Quercia. 2022. Predicting Meeting Success With Nuanced Emotions. *IEEE Pervasive Computing* 21, 2 (April 2022), 51–59. <https://doi.org/10.1109/MPRV.2022.3145047>

APPENDIX

A GLM TABLES FOR SIDE MODELS

The tables shared in this section show detailed effects and interaction effects of a subset of attributes on effectiveness and inclusiveness. These subsets are suggested by the sub-graphs that emerged from graph modeling fitted on all attributes.

	Coef	Standardized Coef	p-value
Intercept	3.80	37.79	0.00
Short Call (30min or less)	-0.27	-2.39	0.02
Meeting Size	-0.06	-8.63	0.00
Short Call (30min or less) : Meeting Size	0.00	-0.50	0.62
Recurring	-0.37	-3.24	0.00
Recurring : Meeting Size	0.03	3.41	0.00

Table 7. GLM results in modeling the probability of Effective by Participation

	Coef	Standardized Coef	p-value
Intercept	-0.40	-6.90	0.00
Meeting Size (8 or less)	2.00	2.00	0.02
VideoDuration30%	0.16	1.84	0.06
CallDuration	0.00	4.22	0.00
Meeting Size (8 or less) : VideoDuration30%	0.46	5.52	0.00
VideoDuration30% : CallDuration	0.00	-2.05	0.04
Meeting Size (8 or less) : VideoDuration30%	0.00	-2.00	0.05

Table 8. GLM results in modeling the probability of Participation by Video usage

Received January 2023; revised July 2023; accepted November 2023

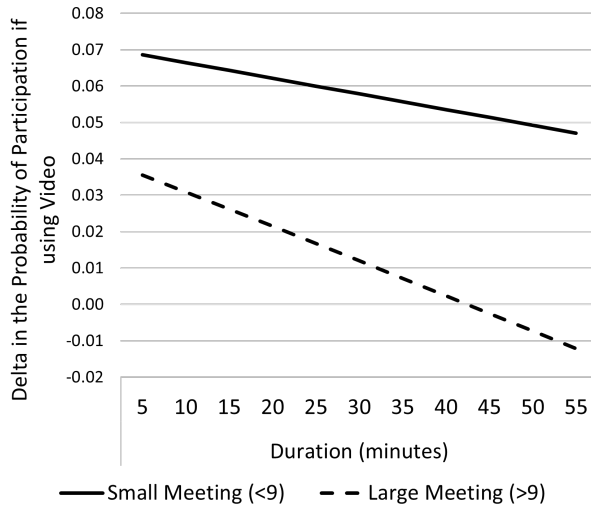


Fig. 9. Interaction effect between Call Duration, Meeting Size, Video usage, and Participation.

Name	Type
Microphone Failure (Initialization)	bool
Microphone Failure (Mid-Call)	bool
Media Failure	bool
Media Ever Flowed	bool
Reconnect Failure	bool
Reconnect Issue	bool
Call Dropped	bool
Video Duration Percent	float
Audio Only	bool
Video Only	bool
Video or ScreenShare	bool
ScreenShare Only	bool
Audio Participation Rate (based on audio packet decoding counts)	float
Is Rater The Meeting Host	bool
Is Friday	bool
is Monday	bool
Is Country GroupA	bool
Is Country GroupB	bool
Meeting Size	int
Call Duration	float
Predicted Probability of Call Quality Issues	float
Total Time In Meeting In The Same Day	float
Total Calls In The Same Day	int
ScreenShare > 10%	bool
Quality Issues	bool
Reliability Issues	bool
Participation	int
Video Duration > 30%	bool
Recurring	bool
ScreenShare	bool
Small Meeting (8 or less)	bool
Short Call (10min. or less)	bool
Long Call (1hr or more)	bool
Headset	bool
Busy Day (10 or More Calls)	bool
Short Hours in Meetings (Less Than 1hr In Calls On The Same Day)	bool

Table 9. List of variables used for predictive modeling.