

APPENDIX

Description of Intervention

Players take on the persona of Andy Jordan, a young hospitalist who moves home after the disappearance of his estranged grandfather, Robert Jordan, and begins a job at a local community hospital. The player has two objectives: to diagnose and treat patients admitted to the hospital, and to solve the mystery of Robert's disappearance.

Patient cases fall into two categories, 'teaching' and 'non-teaching.' Interactions with the 'teaching' patients are designed to communicate a didactic principle that instantiates the game objective of encouraging players to have ACP conversations with all patients over the age of 65 (see **Box**). These patients have a serious illness but are not at the very end-of-life. When players fail to engage in ACP conversations, the patient returns with complications that require additional treatment. Players also receive feedback on their performance from in-game characters (e.g. peers, family members, or their supervisor). The feedback includes factual information about the probability of poor outcomes among patients over 65 who require hospitalization and a reminder about the value of early ACP conversations. In contrast, when players engage in ACP conversations, they subsequently receive an update about the patient's condition, describing how that ACP improved the care of the patient downstream, and a compliment on their decision-making and communication skills. Relevant patients also provide an opportunity for players to observe best practice principles of a high-quality serious illness conversation modeled on Ariadne Lab's Serious Illness Conversation Guide.²⁶ Specifically, when players choose to engage in ACP conversations, the interaction unfolds with Andy asking key questions from the guide and following other best practices (e.g. Andy Jordan pulls up a chair and sits for the conversation).

'Non-teaching' patients either have a critical, immediately life-threatening illness or a diagnostically challenging problem. These cases were designed to increase challenge levels

and associated game-play enjoyment. Players do not receive in-game feedback on their treatment of 'non-teaching' patients. Instead, they receive a summary of their performance on all cases at the end of the game that summarizes decisions made on the teaching cases and the accuracy of their diagnoses for the non-teaching cases.

The mystery component of *Hopewell Hospitalist* occurs concurrently with the clinical challenges, and serves to facilitate players' identification with their character and interest in their task. Players must solve Robert's disappearance through interactions with other characters, including patients, and their physical environment. Andy Jordan's background and character are also revealed through these interactions, which are designed to make him and his decisions more appealing and sympathetic.

Statistical Plan

Here we provide additional information about our analytic plan.

Primary Analysis

Let Y_{ijt} denote the binary outcome variable (coded as 1 if an ACP conversation occurred and 0 otherwise) for patient i seen at hospital j at time t ; $Game_{jt}$ a binary variable indicating whether hospital j has received the Game during period t ($Game_{jt} = 1$ if received by hospital j before or during period t and 0 otherwise), x_{ijt} a vector of patient-level covariates, z_j a vector of hospital-level covariates and θ_j a random effect for hospital. The mathematical specification of the statistical model is given as $Y_{ijt}|\theta_j \sim Bernoulli(\pi_{ijt})$, where

$$\text{logit}(\pi_{ijt}) = \log\left(\frac{\pi_{ijt}}{1 - \pi_{ijt}}\right) = \beta_0 + \beta_{1t} + \beta_2 Game_{jt} + \beta_3 x_{ijt} + \beta_4 z_j + \theta_j$$

where $\theta_j \sim Normal(0, \tau^2)$ is the distribution of the hospital-level random effects to account for the fact that the statistical significance of inferences about the effect of the game are likely to be reduced by the clustering of patients in hospitals. The model includes fixed-effects for time-

period, β_{1t} , to allow for an unstructured trend across calendar time, which makes the effect of the game (the primary target of inference) to be estimated net of any time-trend. The key coefficient of interest is β_2 , which captures the structural shift in the outcome of patients who were enrolled in the study when the hospital receives the iPads, net of general trends across time and other covariates. Because this is a cluster-randomized study, there is a risk that the hospitals in each step are not perfectly balanced, despite attempts to balance these during randomization by forming blocks, and that the distributions of patient characteristics of patients treated by a given hospital may vary across time. To mitigate these concerns, we will adjust for judiciously selected patient and hospital covariates that we hypothesize are reasonably likely to be associated with the outcome. We do not plan to adjust for time-varying hospital-level covariates but we will adjust for whether the hospital was in other programs (e.g., the bundled payment care initiative (BPCI) program) that might influence the culture of the hospital towards ACP; an advantage of adjusting for BPCI participation is that we may obtain more precise inferences.

The reason why physician is excluded from the above model is that a patient may receive care from multiple physicians during their hospital stay. This makes it difficult to designate a single physician as being responsible for the patient's care and thus whether or not they receive an ACP conversation. In our primary analysis we hold the hospital as a collective unit as being responsible for the patient and, therefore, exclude any involvement of physician factors or identifiers in relation to the likelihood of the patient having an ACP conversation. However, based on analyses of preliminary data, we anticipate that for 80% of hospitalizations a single physician will dominate the care of the patient. Therefore, in a sensitivity analysis, we will add a physician layer to the above model and perform a physician-level analysis. Where more than one physician treats a patient, we will assign the patient to the discharging physician, as per the practice of the staffing organization. The resulting statistical model will be a three-level model with physician as the second level (between patient and hospital) to allow patients to be

nested within physicians that are in turn nested within hospitals. Because patients are not randomized to physician, we will consider adjusting for physician covariates, emulating some of the secondary analyses described below.

Secondary analyses

In secondary analyses, we will also explore whether there is evidence on an interaction effect between BPCI participation and the impact of the game on the adjusted odds that a patient has an ACP billed. We will also estimate the effect of the intervention on ACP practices, using both the chart review and the MiPS measures to estimate the sensitivity and specificity of the different methods of measuring ACP. Finally, we will test the effect of mediators on the effect of the intervention on practice patterns, including the dose of a patient's exposure to the intervention, physicians' self-reported engagement with the intervention, and physicians' prior training. A natural game exposure-dose is the number of physicians, encountered by the patient, who had played the game by the time they cared for the patient. The game-exposure measure will replace the hospital-level indicator of game intervention status as the key predictor in these analyses. In analyses in which a single physician is attributed to the patient, the indicator of whether or not that physician has played the game will become the primary predictor of interest, although we may still include other exposure variables in order to extract the independent effect of each source of exposure.

The above factors are potential mediators of the effect of the game being employed at a hospital on patient outcomes as they are on the causal pathway of the hospital-level intervention to patient outcomes; if no physicians who indicated their willingness to participate in the study end up playing the game it is difficult to imagine how the game could then impact their patients' outcomes. Likewise, the hypothesis that a patient who encounters multiple physicians who played the game will have outcomes that are more pronounced than a patient who encountered

only a single physician or even no physicians who played the game a priori appears to be plausible.

In a potential extended analysis, we will adapt statistical methods for incorporating the sensitivity and specificity of the measurement of the occurrence of an ACP conversation, which is informed by the agreement between chart-review and insurance-claim (or MiPS) measurement, into the analysis. The resulting analysis can be viewed as a calibration analysis that combines the standard cluster-randomized stepped-wedge design with a bivariate outcome (a more expensive measurement in the form of chart-review and a less expensive measurement in the form of insurance-claim or MiPS) in order to evaluate the impact of the deployment of the game at a hospital on chart-based measurement of ACP occurrence. The statistical model entwining the outcomes will allow the missing values of chart-based measurement for those observations where charts are not reviewed to be learned from observations for which multiple forms of ACP measurement are made and automatically allow for uncertainty in the missing values of chart-review measurements to permeate through the analysis. A Bayesian statistical model and Bayesian computational methods may provide the least burdensome pathway to successfully implementing this analysis.

Power calculation

We arrived at our sample size using a combination of feasibility (cost) and assumptions regarding effect size, absent any pilot data about the latter. For each step, we plan to recruit between 25 to 30 physicians from each of 4 to 8 hospitals. Assuming a baseline ACP rate of 22% (rising by 1.5 percentage-points per-quarter), a hospital intra-class correlation (ICC) coefficient of 0.01-0.10, and 160 evaluable patients per physician-quarter, we can detect a 3.5 percentage-point difference between ACP practices before and after the distribution of the intervention using a two-sided test at the 0.05-level with power in excess of 99%, even under the most conservative sample-size assumptions. If we invert the problem to find the smallest

effect-size at which our study has 80% power, we find that in the most conservative scenario (76,800 total patients) we can detect a 1.5 percentage-point difference and in the most optimistic scenario (192,000 total patients), we can detect a 1 percentage-point increase.

The method of computing power for this stepped-wedge design follows the commonly used strategy for cluster randomized trials of first determining the design-effect, which can be thought of as a measure of the inefficiency of the given design in comparison to a completely randomized design that is expressed in terms of a ratio of the sample-sizes needed to obtain equally precise estimates, and then applying conventional power calculations. The latter computes power for a two-population comparison using the effective-sample-sizes determined from the design-effect. We estimate the design-effect using the expression in Woertman et al (2013), that was clarified and illustrated in Hemming (2016). Because hospitals may induce correlations in the outcomes of patients who receive care from them, we perform illustrative power calculations that account for the net impact of clustering at the hospital-level. Based on our own prior research and published results of others, we decided that the ICC of hospital is highly likely to be in the range 0.01 to 0.10. The design-effects across the optimistic and pessimistic scenarios ranged between 2.88 and 3.14, implying that for all considered scenarios the stepped-wedge design is about 33% as efficient as a patient-level completely randomized design. The effective sample-sizes (ESS) per group ranged from 30,603 to 12,388 patients per group over the study period (the 5 steps and a baseline period).

The second part of the calculation is to determine the power of a two-group comparison of a binary outcome in the absence of clustering when the total sample-size per group equals the above values for the ESS. Because the sample-sizes are still reasonably large, an asymptotic normal approximation is well justified, especially at a baseline ACP rate of 22%. Because we generally err on the side of making conservative estimates about the level of information available (e.g., we may extend the baseline period in which can retrospectively

acquire data to 3-months), Therefore, this approximate two-step calculation yields trustworthy estimates of power that, if anything, are expected to err on the side of being conservative.