

¹ Text embedding models yield high-resolution insights
² into conceptual knowledge from short multiple-choice
³ quizzes

⁴ Paxton C. Fitzpatrick¹, Andrew C. Heusser^{1, 2}, and Jeremy R. Manning^{1, *}

¹Department of Psychological and Brain Sciences

Dartmouth College, Hanover, NH 03755, USA

²Akili Interactive Labs

Boston, MA 02110, USA

*Corresponding author: Jeremy.R.Manning@Dartmouth.edu

⁵ **Abstract**

⁶ We develop a mathematical framework, based on natural language processing models, for track-
⁷ ing and characterizing the acquisition of conceptual knowledge. Our approach embeds each
⁸ concept in a high-dimensional representation space, where nearby coordinates reflect similar or
⁹ related concepts. We test our approach using behavioral data from participants who answered
¹⁰ small sets of multiple-choice quiz questions interleaved between watching two course videos
¹¹ from the Khan Academy platform. We apply our framework to the videos' transcripts and
¹² the text of the quiz questions to quantify the content of each moment of video and each quiz
¹³ question. We use these embeddings, along with participants' quiz responses, to track how the
¹⁴ learners' knowledge changed after watching each video. Our findings show how a small set of
¹⁵ quiz questions may be used to obtain rich and meaningful high-resolution insights into what
¹⁶ each learner knows, and how their knowledge changes over time as they learn.

¹⁷ **Keywords:** education, learning, knowledge, concepts, natural language processing

¹⁸ **Introduction**

¹⁹ Suppose that a teacher had access to a complete, tangible “map” of everything a student knows.
²⁰ Defining what such a map might even look like, let alone how it might be constructed or filled in, is
²¹ itself a non-trivial problem. But if a teacher *were* to gain access to such a map, how might it change
²² their ability to teach that student? Perhaps they might start by checking how well the student
²³ knows the to-be-learned information already, or how much they know about related concepts.
²⁴ For some students, they could potentially optimize their teaching efforts to maximize efficiency
²⁵ by focusing primarily on not-yet-known content. For other students (or other content areas), it
²⁶ might be more effective to optimize for direct connections between already known content and
²⁷ new material. Observing how the student’s knowledge changed over time, in response to their
²⁸ teaching, could also help to guide the teacher towards the most effective strategy for that individual
²⁹ student.

³⁰ A common approach to assessing a student’s knowledge is to present them with a set of quiz
³¹ questions, calculate the proportion they answer correctly, and provide them with feedback in the
³² form of a simple numeric or letter grade. While such a grade can provide *some* indication of whether
³³ the student has mastered the to-be-learned material, any univariate measure of performance on a
³⁴ complex task sacrifices certain relevant information, risks conflating underlying factors, and so on.
³⁵ For example, consider the relative utility of the theoretical map described above that characterizes
³⁶ a student’s knowledge in detail, versus a single annotation saying that the student answered 85%
³⁷ of their quiz questions correctly, or that they received a ‘B’. Here, we show that the same quiz data
³⁸ required to compute proportion-correct scores or letter grades can instead be used to obtain far
³⁹ more detailed insights into what a student knew at the time they took the quiz.

⁴⁰ Designing and building procedures and tools for mapping out knowledge touches on deep
⁴¹ questions about what it means to learn. For example, how do we acquire conceptual knowledge?
⁴² Memorizing course lectures or textbook chapters by rote can lead to the superficial *appearance*
⁴³ of understanding the underlying content, but achieving true conceptual understanding seems to
⁴⁴ require something deeper and richer. Does conceptual understanding entail connecting newly

45 acquired information to the scaffolding of one’s existing knowledge or experience [6, 11, 13, 15, 30,
46 65]? Or weaving a lecture’s atomic elements (e.g., its component words) into a structured network
47 that describes how those individual elements are related [40, 70]? Conceptual understanding
48 could also involve building a mental model that transcends the meanings of those individual
49 atomic elements by reflecting the deeper meaning underlying the gestalt whole [37, 41, 62, 69].

50 The difference between “understanding” and “memorizing,” as framed by researchers in ed-
51 ucation, cognitive psychology, and cognitive neuroscience [e.g., 23, 28, 33, 41, 62], has profound
52 analogs in the fields of natural language processing and natural language understanding. For
53 example, considering the raw contents of a document (e.g., its constituent symbols, letters, and
54 words) might provide some clues as to what the document is about, just as memorizing a passage
55 might provide some ability to answer simple questions about it. However, text embedding mod-
56 els [e.g., 7, 8, 10, 12, 16, 39, 50, 71] also attempt to capture the deeper meaning *underlying* those
57 atomic elements. These models consider not only the co-occurrences of those elements within and
58 across documents, but (in many cases) also patterns in how those elements appear across different
59 scales (e.g., sentences, paragraphs, chapters, etc.), the temporal and grammatical properties of the
60 elements, and other high-level characteristics of how they are used [42, 43]. To be clear, this is not
61 to say that text embedding models themselves are capable of “understanding” deep conceptual
62 meaning in any traditional sense. But rather, their ability to capture the underlying *structure* of
63 text documents beyond their surface-level contents provides a computational framework through
64 which those documents’ deeper conceptual meanings may be quantified, explored, and under-
65 stood. According to these models, the deep conceptual meaning of a document may be captured
66 by a feature vector in a high-dimensional representation space, wherein nearby vectors reflect con-
67 ceptually related documents. A model that succeeds at capturing an analogue of “understanding”
68 is able to assign nearby feature vectors to two conceptually related documents, *even when the specific*
69 *words contained in those documents have limited overlap*. In this way, “concepts” are defined implicitly
70 by the model’s geometry [e.g., how the embedding coordinate of a given word or document relates
71 to the coordinates of other text embeddings; 56].

72 Given these insights, what form might a representation of the sum total of a person’s knowledge

73 take? First, we might require a means of systematically describing or representing (at least some
74 subset of) the nearly infinite set of possible things a person could know. Second, we might want to
75 account for potential associations between different concepts. For example, the concepts of “fish”
76 and “water” might be associated in the sense that fish live in water. Third, knowledge may have
77 a critical dependency structure, such that knowing about a particular concept might require first
78 knowing about a set of other concepts. For example, understanding the concept of a fish swimming
79 in water first requires understanding what fish and water *are*. Fourth, as we learn, our “current
80 state of knowledge” should change accordingly. Learning new concepts should both update our
81 characterizations of “what is known” and also unlock any now-satisfied dependencies of those
82 newly learned concepts so that they are “tagged” as available for future learning.

83 Here we develop a framework for modeling how conceptual knowledge is acquired during
84 learning. The central idea behind our framework is to use text embedding models to define the
85 coordinate systems of two maps: a *knowledge map* that describes the extent to which each concept is
86 currently known, and a *learning map* that describes changes in knowledge over time. Each location
87 on these maps represents a single concept, and the maps’ geometries are defined such that related
88 concepts are located nearby in space. We use this framework to analyze and interpret behavioral
89 data collected from an experiment that had participants answer sets of multiple-choice questions
90 about a series of recorded course lectures.

91 Our primary research goal is to advance our understanding of what it means to acquire deep,
92 real-world conceptual knowledge. Traditional laboratory approaches to studying learning and
93 memory (e.g., list-learning studies) often draw little distinction between memorization and under-
94 standing. Instead, these studies typically focus on whether information is effectively encoded or
95 retrieved, rather than whether the information is *understood*. Approaches to studying conceptual
96 learning, such as category learning experiments, can begin to investigate the distinction between
97 memorization and understanding, often by training participants to distinguish arbitrary or random
98 features in otherwise meaningless categorized stimuli [1, 20, 21, 24, 31, 59]. However the objective
99 of real-world training, or learning from life experiences more generally, is often to develop new
100 knowledge that may be applied in *useful* ways in the future. In this sense, the gap between modern

learning theories and modern pedagogical approaches that inform classroom learning strategies is enormous: most of our theories about *how* people learn are inspired by experimental paradigms and models that have only peripheral relevance to the kinds of learning that students and teachers actually seek [28, 41]. To help bridge this gap, our study uses course materials from real online courses to inform, fit, and test models of real-world conceptual learning. We also provide a demonstration of how our models can be used to construct “maps” of what students know, and how their knowledge changes with training. In addition to helping to visually capture knowledge (and changes in knowledge), we hope that such maps might lead to real-world tools for improving how we educate. Taken together, our work shows that existing course materials and evaluative tools like short multiple-choice quizzes may be leveraged to gain highly detailed insights into what students know and how they learn.

Results

At its core, our main modeling approach is based around a simple assumption that we sought to test empirically: all else being equal, knowledge about a given concept is predictive of knowledge about similar or related concepts. From a geometric perspective, this assumption implies that knowledge is fundamentally “smooth.” In other words, as one moves through a space representing an individual’s knowledge (where similar concepts occupy nearby coordinates), their “level of knowledge” should change relatively gradually. To begin to test this smoothness assumption, we sought to track participants’ knowledge and how it changed over time in response to training. Two overarching goals guide our approach. First, we want to gain detailed insights into what learners know at different points in their training. For example, rather than simply reporting on the proportions of questions participants answer correctly (i.e., their overall performance), we seek estimates of their knowledge about a variety of specific concepts. Second, we want our approach to be potentially scalable to large numbers of diverse concepts, courses, and students. This requires that the conceptual content of interest be discovered *automatically*, rather than relying on manually produced ratings or labels.



Figure 1: Experimental paradigm. Participants alternate between completing three 13-question multiple-choice quizzes and watching two Khan Academy lectures. Each quiz contains a mix of 5 questions about Lecture 1, 5 questions about Lecture 2, and 3 questions about general physics knowledge. The specific questions reflected on each quiz, and the orders of each quiz's questions, were randomized across participants.

127 We asked participants in our study to complete brief multiple-choice quizzes before, between,
 128 and after watching two lecture videos from the Khan Academy [36] platform (Fig. 1). The first
 129 lecture video, entitled *Four Fundamental Forces*, discussed the four fundamental forces in physics:
 130 gravity, strong and weak interactions, and electromagnetism. The second, entitled *Birth of Stars*,
 131 provided an overview of our current understanding of how stars form. We selected these particular
 132 lectures to satisfy three general criteria. First, we wanted both lectures to be accessible to a broad
 133 audience (i.e., with minimal prerequisite knowledge) so as to limit the impact of prior training
 134 on participants' abilities to learn from the lectures. To this end, we selected two introductory
 135 videos that were intended to be viewed at the start of students' training in their respective content
 136 areas. Second, we wanted the two lectures to have some related content, so that we could test
 137 our approach's ability to distinguish similar conceptual content. To this end, we chose two videos
 138 from the same Khan Academy course domain, "Cosmology and Astronomy." Third, we sought to
 139 minimize dependencies and specific overlap between the videos. For example, we did not want
 140 participants' abilities to understand one video to (directly) influence their abilities to understand the
 141 other. To satisfy this last criterion, we chose videos from two different lecture series (Lectures 1 and
 142 2 were from the "Scale of the Universe" and "Stars, Black Holes, and Galaxies" series, respectively).

143 We also wrote a set of multiple-choice quiz questions that we hoped would enable us to
 144 evaluate participants' knowledge about each individual lecture, along with related knowledge



Figure 2: Modeling course content. **A. Building a document pool from sliding windows of text.** We decompose each lecture’s transcript into a series of overlapping sliding windows. The full set of transcript snippets (across all windows) may be treated as a set of “documents” for training a text embedding model. **B. Constructing lecture content trajectories.** After training the model on the sliding windows from both lectures, we transform each lecture into a “trajectory” through text embedding space by joining the embedding coordinates of successive sliding windows parsed from its transcript. **C. Embedding multiple lectures and questions in a shared space.** We apply the same model (trained on the two lectures’ windows) to both lectures, along with the text of each question in our pool (Supp. Tab. 1), to project them into a shared text embedding space. This results in one trajectory per lecture and one coordinate for each question. Here, we have projected the 15-dimensional embeddings onto their first 3 principal components for visualization.

145 about physics concepts not specifically presented in either video (see Supp. Tab. 1 for the full list
 146 of questions in our stimulus pool). Participants answered questions randomly drawn from each
 147 content area (Lecture 1, Lecture 2, and general physics knowledge) on each of the three quizzes.
 148 Quiz 1 was intended to assess participants’ “baseline” knowledge before training, Quiz 2 assessed
 149 knowledge after watching the *Four Fundamental Forces* video (i.e., Lecture 1), and Quiz 3 assessed
 150 knowledge after watching the *Birth of Stars* video (i.e., Lecture 2).

151 To study in detail how participants’ conceptual knowledge changed over the course of the
 152 experiment, we first sought to model the conceptual content presented to them at each moment
 153 throughout each of the two lectures. We adapted an approach we developed in prior work [29]
 154 to identify the latent themes in the lectures using a topic model [8]. Briefly, topic models take
 155 as input a collection of text documents, and learn a set of “topics” (i.e., latent themes) from their
 156 contents. Once fit, a topic model can be used to transform arbitrary (potentially new) documents
 157 into sets of “topic proportions,” describing the weighted blend of learned topics reflected in their
 158 texts. We parsed automatically generated transcripts of the two lectures into overlapping sliding
 159 windows, where each window contained the text of the lecture transcript from a particular time

span. We treated the set of text snippets (across all of these windows) as documents to fit the model (Fig. 2A; see *Constructing text embeddings of multiple lectures and questions*). Transforming the text from every sliding window with the model yielded a number-of-windows by number-of-topics (15) topic-proportions matrix describing the unique mixture of broad themes from both lectures reflected in each window’s text. Each window’s “topic vector” (i.e., column of the topic-proportions matrix) is analogous to a coordinate in a 15-dimensional space whose axes are topics discovered by the model. Within this space, each lecture’s sequence of topic vectors (i.e., corresponding to its transcript’s overlapping text snippets across sliding windows) forms a *trajectory* that captures how its conceptual content unfolds over time (Fig. 2B). We resampled these trajectories to a resolution of one topic vector for each second of video (i.e., 1 Hz).

We hypothesized that a topic model trained on transcripts of the two lectures should also capture the conceptual knowledge probed by each quiz question. If indeed the topic model could capture information about the deeper conceptual content of the lectures (i.e., beyond surface-level details such as particular word choices), then we should be able to recover a correspondence between each lecture and questions *about* each lecture. Importantly, such a correspondence could not solely arise from superficial text matching between lecture transcripts and questions, since the lectures and questions often used different words (Supp. Fig. 5) and phrasings. Simply comparing the average topic weights from each lecture and question set (averaging across time and questions, respectively) reveals a striking correspondence (Supp. Fig. 2). Specifically, the average topic weights from Lecture 1 are strongly correlated with the average topic weights from Lecture 1 questions ($r(13) = 0.809$, $p < 0.001$, 95% confidence interval (CI) = [0.633, 0.962]), and the average topic weights from Lecture 2 are strongly correlated with the average topic weights from Lecture 2 questions ($r(13) = 0.728$, $p = 0.002$, 95% CI = [0.456, 0.920]). At the same time, the average topic weights from the two lectures are *negatively* correlated with the average topic weights from their non-matching question sets (Lecture 1 video vs. Lecture 2 questions: $r(13) = -0.547$, $p = 0.035$, 95% CI = [-0.812, -0.231]; Lecture 2 video vs. Lecture 1 questions: $r(13) = -0.612$, $p = 0.015$, 95% CI = [-0.874, -0.281]), indicating that the topic model also exhibits some degree of specificity. The full set of pairwise comparisons between average topic weights for the lectures and question sets



Figure 3: Lecture and question topic overlap. A. Topic weight variability. The bar plots display the variance of each topic's weight across lecture timepoints (top row) and questions (bottom row); colors denote topics. The top-weighted words from the most “expressive” (i.e., variable across observations) topic from each lecture are displayed in the upper right (orange: topic 2; yellow-green: topic 4). The top-weighted words from the full set of topics may be found in Supplementary Table 2. **B. Relationships between topic weight variability.** Pairwise correlations between the distributions of topic weight variance for each lecture and question set. Each row and column corresponds to a bar plot in Panel A.

188 is reported in Supplementary Figure 2.

189 Another, more sensitive, way of summarizing the conceptual content of the lectures and questions is to look at *variability* in how topics are weighted over time and across different questions 190 (Fig. 3). Intuitively, the variability in the expression of a given topic relates to how much “information” [22] the lecture (or question set) reflects about that topic. For example, suppose a given 191 topic is weighted on heavily throughout a lecture. That topic might be characteristic of some 192 aspect or property of the lecture *overall* (conceptual or otherwise), but unless the topic’s weights 193 changed in meaningful ways over time, the topic would be a poor indicator of any *specific* conceptual 194 content in the lecture. We therefore also compared the variances in topic weights (across time 195 or questions) between the lectures and questions. The variability in topic expression (over time 196 or questions) was similar for the Lecture 1 video and questions ($r(13) = 0.824, p < 0.001$, 197 95% CI = [0.696, 0.973]) and the Lecture 2 video and questions ($r(13) = 0.801, p < 0.001$, 95% 198 CI = [0.539, 0.958]). Simultaneously, as reported in Figure 3B, the variabilities in topic expression 199 across *different* videos and lecture-specific questions (i.e., Lecture 1 video vs. Lecture 2 questions; 200 201

202 Lecture 2 video vs. Lecture 1 questions) were negatively correlated, and neither video’s topic
203 variability was reliably correlated with the topic variability across general physics knowledge
204 questions. Taken together, the analyses reported in Figure 3 and Supplementary Figure 2 indicate
205 that a topic model fit to the videos’ transcripts can also reveal correspondences (at a coarse scale)
206 between the lectures and questions.

207 While an individual lecture may be organized around a single broad theme at a coarse scale,
208 at a finer scale, each moment of a lecture typically covers a narrower range of content. Given the
209 correspondence we found between the variabilities in topic expression across moments of each
210 lecture and questions from its corresponding set (Fig. 3), we wondered whether the text embedding
211 model might additionally capture these conceptual relationships at a finer scale. For example, if a
212 particular question asks about the content from one small part of a lecture, we wondered whether
213 the text embeddings could be used to automatically identify the “matching” moment(s) in the
214 lecture. To explore this, we computed the correlation between each question’s topic weights
215 and the topic weights for each second of its corresponding lecture, and found that each question
216 appeared to be temporally specific (Fig. 4). In particular, most questions’ topic vectors were
217 maximally correlated with a well-defined (and relatively narrow) range of timepoints from their
218 corresponding lectures, and the correlations fell off sharply outside of that range (Supp. Figs. 3, 4).
219 We also qualitatively examined the best-matching intervals for each question by comparing the
220 question’s text to the transcribed text from the most-correlated parts of the lectures (Supp. Tab. 3).
221 Despite that the questions were excluded from the text embedding model’s training set, in general
222 we found (through manual inspection) a close correspondence between the conceptual content
223 that each question probed and the content covered by the best-matching moments of the lectures.
224 Two representative examples are shown at the bottom of Figure 4.

225 The ability to quantify how much each question is “asking about” the content from each moment
226 of the lectures could enable high-resolution insights into participants’ knowledge. Traditional
227 approaches to estimating how much a student “knows” about the content of a given lecture entail
228 administering some form of assessment (e.g., a quiz) and computing the proportion of correctly
229 answered questions. But if two students receive identical scores on such an exam, might our



Figure 4: Which parts of each lecture are captured by each question? Each panel displays time series plots showing how each question’s topic vector correlates with each video timepoint’s topic vector (Panel A.: correlations for the *Four Fundamental Forces* lecture and associated questions; Panel B.: correlations for the *Birth of Stars* lecture and associated questions). The colors denote question identities. The diamonds in each panel denote the moment of peak correlation between the indicated question and the lecture trajectory. The associated questions’ text and snippets of the lectures’ transcripts from the surrounding 30 seconds, are displayed at the bottom of the figure.

modeling framework help us to gain more nuanced insights into the *specific* content that each student has mastered (or failed to master)? For example, a student who misses three questions that were all about the same concept (e.g., concept *A*) will have gotten the same *proportion* of questions correct as another student who missed three questions about three *different* concepts (e.g., *A*, *B*, and *C*). But if we wanted to help these two students fill in the “gaps” in their understandings, we might do well to focus specifically on concept *A* for the first student, but to also add in materials pertaining to concepts *B* and *C* for the second student. In other words, raw “proportion-correct” measures may capture *how much* a student knows, but not *what* they know. We wondered whether our modeling framework might enable us to (formally and automatically) infer participants’ knowledge at the scale of individual concepts (e.g., as captured by a single moment of a lecture).

We developed a simple formula (Eqn. 1) for using a participant’s responses to a small set of multiple-choice questions to estimate how much the participant “knows” about the concept reflected by any arbitrary coordinate x in text embedding space (e.g., the content reflected by

any moment in a lecture they had watched; see *Estimating dynamic knowledge traces*). Essentially, the estimated knowledge at coordinate x is given by the weighted proportion of quiz questions the participant answered correctly, where the weights reflect how much each question is “about” the content at x . When we apply this approach to estimate the participant’s knowledge about the content presented in each moment of each lecture, we can obtain a detailed time course describing how much “knowledge” that participant has about the content presented at any part of the lecture. As shown in Figure 5A and C, we can apply this approach separately for the questions from each quiz participants took throughout the experiment. From just a few questions per quiz (see *Estimating dynamic knowledge traces*), we obtain a high-resolution snapshot (at the time each quiz was taken) of what the participants knew about any moment’s content, from either of the two lectures they watched (comprising a total of 1,100 samples across the two lectures).

While the time courses in Figure 5A and C provide detailed *estimates* about participants’ knowledge, these estimates are of course only *useful* to the extent that they accurately reflect what participants actually know. As one sanity check, we anticipated that the knowledge estimates should reflect a content-specific “boost” in participants’ knowledge after watching each lecture. In other words, if participants learn about each lecture’s content upon watching it, the knowledge estimates should capture that. After watching the *Four Fundamental Forces* lecture, participants should exhibit more knowledge for the content of that lecture than they had before, and that knowledge should persist for the remainder of the experiment. Specifically, knowledge about that lecture’s content should be relatively low when estimated using Quiz 1 responses, but should increase when estimated using Quiz 2 or 3 responses (Fig. 5B). Indeed, we found that participants’ estimated knowledge about the content of *Four Fundamental Forces* was substantially higher on Quiz 2 versus Quiz 1 ($t(49) = 8.764, p < 0.001$) and on Quiz 3 versus Quiz 1 ($t(49) = 10.519, p < 0.001$). We found no reliable differences in estimated knowledge about that lecture’s content on Quiz 2 versus 3 ($t(49) = 0.160, p = 0.874$). Similarly, we hypothesized (and subsequently confirmed) that participants should show greater estimated knowledge about the content of the *Birth of Stars* lecture after (versus before) watching it (Fig. 5D). Specifically, since participants watched that lecture after taking Quiz 2 (but before Quiz 3), we hypothesized that their knowledge estimates



Figure 5: Estimating knowledge about the content presented at each moment of each lecture. **A. Knowledge about the time-varying content of *Four Fundamental Forces*.** Each trace displays the weighted proportion of correctly answered questions about the content reflected in each moment of the lecture (see *Estimating dynamic knowledge traces*), using responses from a single quiz (color). The traces are averaged across participants. **B. Average estimated knowledge about *Four Fundamental Forces*.** Each bar displays the across-timepoint average knowledge, estimated using the responses to one quiz's questions. **C. Knowledge about the time-varying content of *Birth of Stars*.** The panel is in the same format as Panel A, but here the knowledge estimates are for the moment-by-moment content of the *Birth of Stars* lecture. **D. Average estimated knowledge about *Birth of Stars*.** The panel is in the same format as Panel B, but here the knowledge estimates are for the content of the *Birth of Stars* lecture. All panels: error ribbons and error bars denote 95% confidence intervals, estimated across participants.

271 should be relatively low on Quizzes 1 and 2, but should show a “boost” on Quiz 3. Consistent
272 with this prediction, we found no reliable differences in estimated knowledge about the *Birth of*
273 *Stars* lecture content on Quizzes 1 versus 2 ($t(49) = 1.013, p = 0.316$), but the estimated knowl-
274 edge was substantially higher on Quiz 3 versus 2 ($t(49) = 10.561, p < 0.001$) and Quiz 3 versus 1
275 ($t(49) = 8.969, p < 0.001$).

276 If we are able to accurately estimate a participant’s knowledge about the content tested by a
277 given question, our estimates of their knowledge should carry some predictive information about
278 whether they are likely to answer that question correctly or incorrectly. We developed a statistical
279 approach to test this claim. For each quiz question a participant answered, in turn, we used
280 Equation 1 to estimate their knowledge at the given question’s embedding space coordinate based
281 on other questions that participant answered on the same quiz. We repeated this for all participants,
282 and for each of the three quizzes. Then, separately for each quiz, we fit a generalized linear mixed
283 model (GLMM) with a logistic link function to explain the likelihood of correctly answering a
284 question as a function of estimated knowledge for its embedding coordinate, while accounting
285 for random variation among participants and questions (see *Generalized linear mixed models*). To
286 assess the predictive value of the knowledge estimates, we compared each GLMM to an analogous
287 (i.e., nested) “null” model that did not consider estimated knowledge using parametric bootstrap
288 likelihood-ratio tests.

289 We carried out three different versions of the analyses described above, wherein we considered
290 different sources of information in our estimates of participants’ knowledge for each quiz question.
291 First, we estimated knowledge at each question’s embedding coordinate using *all* other questions
292 answered by the same participant on the same quiz (“All questions”; Fig. 6, top row). This test was
293 intended to assess the overall predictive power of our approach. Second, we estimated knowledge
294 for each question about a given lecture using only the other questions (from the same participant
295 and quiz) about that *same* lecture (“Within-lecture”; Fig. 6, middle rows). This test was intended to
296 assess the *specificity* of our approach by asking whether our predictions could distinguish between
297 questions about different content covered by the same lecture. Third, we estimated knowledge
298 for each question about one lecture using only questions (from the same participant and quiz)

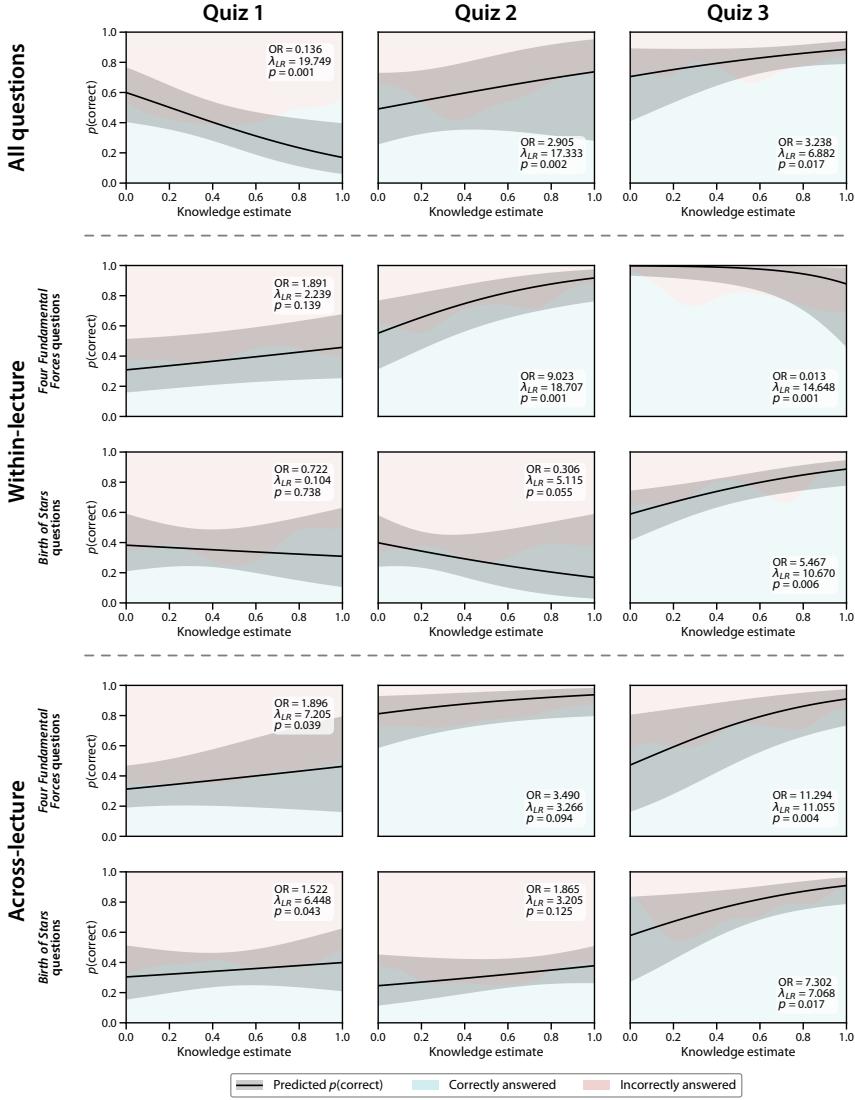


Figure 6: Predicting success on held-out questions using estimated knowledge. We used generalized linear mixed models (GLMMs) to model the likelihood of correctly answering a quiz question as a function of estimated knowledge for its embedding coordinate (see *Generalized linear mixed models*). Separately for each quiz (column), we examined this relationship based on three different sets of knowledge estimates: knowledge for each question based on all other questions the same participant answered on the same quiz (“All questions”; top row), knowledge for each question about one lecture based on all other questions (from the same participant and quiz) about the *same* lecture (“Within-lecture”; middle rows), and knowledge for each question about one lecture based on all questions (from the same participant and quiz) about the *other* lecture (“Across-lecture”; bottom rows). The backgrounds in each panel display kernel density estimates of the relative observed proportions of correctly (blue) versus incorrectly (red) answered questions, for each level of estimated knowledge along the x -axis. The black curves display the (population-level) GLMM-predicted probabilities of correctly answering a question as a function of estimated knowledge. Error ribbons denote 95% confidence intervals.

299 about the *other* lecture (“Across-lecture”; Fig. 6, bottom rows). This test was intended to assess the
300 *generalizability* of our approach by asking whether our predictions held across the content areas of
301 the two lectures.

302 In performing these analyses, our null hypothesis is that the knowledge estimates we compute
303 based on the quiz questions’ embedding coordinates do *not* provide useful information about
304 participants’ abilities to answer those questions. What result might we expect to see if this is the
305 case? To gain an intuition for this scenario, consider the expected outcome if we carried out these
306 same analyses using a simple proportion-correct measure in lieu of our knowledge estimates.
307 Suppose a participant correctly answered n out of q questions on a given quiz. If we hold out
308 a single *correctly* answered question, the proportion of remaining questions answered correctly
309 would be $\frac{n-1}{q-1}$. Whereas if we hold out a single *incorrectly* answered question, the proportion of
310 remaining questions answered correctly would be $\frac{n}{q-1}$. In this way, the proportion of correctly
311 answered remaining questions is always *lower* when the held-out question was answered correctly
312 than when it was answered incorrectly. Because our knowledge estimates are computed as a
313 weighted version of this same proportion-correct score (where each held-in question’s weight
314 reflects its embedding-space distance from the held-out question; see Eqn. 1), if these weights
315 are uninformative (e.g., randomly distributed), then we should expect to see this same inverse
316 relationship between estimated knowledge and performance, on average. On the other hand,
317 if the spatial relationships among the quiz questions’ embeddings *are* predictive of participants’
318 knowledge about the questions’ content, then we would expect *higher* estimated knowledge for
319 held-out correctly (versus incorrectly) answered questions.

320 Before presenting our results, it is worth considering three possible explanations of why a
321 participant might answer a given question correctly or incorrectly. One possibility is that the
322 participant simply *guessed* the answer. A second is that they selected an answer by mistake, despite
323 “knowing” the correct answer. In both of these scenarios, the participant’s knowledge about the
324 question’s content should be uninformative about their observed response. A third possibility is
325 that the participant’s response reflects their *actual* knowledge about the question’s content. In this
326 case, we *might* expect to see a positive relationship between the participant’s knowledge and their

327 likelihood of answering the question correctly. However, in order to see this positive relationship,
328 the participant's knowledge must be structured in a way that is reflected (at least partially) by the
329 embedding space. In other words, if the participant's performance reflects their true knowledge,
330 but our text embedding space does not sufficiently capture the structure of that knowledge, then
331 the knowledge estimates we generate will not be predictive of the participant's performance. In the
332 extreme, if the embedding space is completely unstructured with respect to the content of the quiz
333 questions, then we would expect to see the negative relationship between estimated knowledge
334 and performance that we described above.

335 When we fit a GLMM to estimates of participants' knowledge for each Quiz 1 question based
336 on all other Quiz 1 questions, we observed an outcome consistent with our null hypothesis: higher
337 estimated knowledge at the embedding coordinate of a held-out question was associated with a
338 lower likelihood of answering the question correctly (odds ratio (OR) = 0.136, likelihood-ratio test
339 statistic (λ_{LR}) = 19.749, 95% CI = [14.352, 26.545], p = 0.001). This outcome suggests that our
340 knowledge estimates do *not* provide useful information about participants' Quiz 1 performance
341 when we aggregated across all question content areas. We speculated that this might either
342 indicate that the knowledge estimates are uninformative in general, or about Quiz 1 performance
343 in particular. This would be expected, for example, if participants were guessing about the answers
344 to the Quiz 1 questions (prior to having watched either lecture). When we repeated this analysis
345 for Quizzes 2 and 3, we found that *higher* estimated knowledge for a given question predicted
346 a greater likelihood of answering it correctly (Quiz 2: OR = 2.905, λ_{LR} = 17.333, 95% CI =
347 [14.966, 29.309], p = 0.002; Quiz 3: OR = 3.238, λ_{LR} = 6.882, 95% CI = [6.228, 8.184], p = 0.017).
348 Taken together, these results suggest that our knowledge estimates reliably predict participants'
349 performance on individual held-out quiz questions, but only after participants have received at
350 least some training.

351 We observed a similar pattern of results when used this approach to estimate participants'
352 knowledge about held-out questions from one lecture using their performance on other questions
353 from the *same* lecture. Specifically, for Quiz 1 questions (i.e., prior to watching either), participants'
354 estimated knowledge for the embedding coordinates of held-out *Four Fundamental Forces*-related

355 questions estimated using other *Four Fundamental Forces*-related questions did not reliably pre-
356 dict whether those questions were answered correctly ($OR = 1.891$, $\lambda_{LR} = 2.293$, 95% CI =
357 [2.091, 2.622], $p = 0.139$). The same was true of knowledge estimates for held-out *Birth of Stars*-
358 related questions based on other *Birth of Stars*-related questions ($OR = 0.722$, $\lambda_{LR} = 5.115$, 95% CI =
359 [0.094, 0.146], $p = 0.738$). As in our analysis that included all questions, we speculate that
360 these “null” results might reflect some degree of random guessing on Quiz 1. When we re-
361 peated these within-lecture analyses using questions from Quiz 2 (which participants took im-
362 mediately after viewing *Four Fundamental Forces* but prior to viewing *Birth of Stars*), we found
363 that they now reliably predicted success on *Four Fundamental Forces*-related questions ($OR =$
364 9.023, $\lambda_{LR} = 18.707$, 95% CI = [10.877, 22.222], $p = 0.001$) but not on *Birth of Stars*-related
365 questions ($OR = 0.306$, $\lambda_{LR} = 5.115$, 95% CI = [4.624, 5.655], $p = 0.055$). Here, we speculate
366 that participants might have been guessing about the *Birth of Stars* content (e.g., prior to having
367 watched it), whereas they might have been drawing on some structured knowledge about the *Four*
368 *Fundamental Forces* content (e.g., from having just watched it). When we applied this approach to
369 Quiz 3 responses (given immediately after viewing *Birth of Stars*), we found that within-lecture
370 knowledge estimates for *Birth of Stars*-related questions could now reliably predict success on those
371 questions ($OR = 5.467$, $\lambda_{LR} = 10.670$, 95% CI = [7.998, 12.532], $p = 0.006$). However, within-lecture
372 knowledge estimates for *Four Fundamental Forces* questions answered on Quiz 3 were no longer
373 directly related to the likelihood of successfully answering them and instead exhibited the inverse
374 relationship we would expect to arise from unstructured knowledge (with respect to the embed-
375 ding space; $OR = 0.013$, $\lambda_{LR} = 14.648$, 95% CI = [10.695, 23.096], $p = 0.001$). Speculatively, we
376 suggest that this may reflect participants forgetting some of the *Four Fundamental Forces* content
377 (e.g., perhaps in favor of prioritizing encoding the just-watched *Birth of Stars* content in preparation
378 for the third quiz). If this forgetting happens in a relatively “random” way (with respect to spatial
379 distance within the embedding space), then it could explain why some held-out questions about
380 *Four Fundamental Forces* were answered incorrectly, even if questions at nearby coordinates (i.e.,
381 about similar content) were answered correctly. This might lead our approach to over-estimate
382 knowledge for held-out questions about “forgotten” knowledge that participants answered in-

383 correctly. Taken together, these within-lecture results suggest that our approach can distinguish
384 between questions about different content covered by a single lecture when participants have suf-
385 ficiently structured knowledge about its contents, though this specificity may decrease with time
386 since the relevant material was learned.

387 Finally, we used this approach to estimate participants' knowledge about held-out questions
388 from one lecture using their performance on questions from the *other* lecture. Here we again
389 observed a similar pattern of results, though with some notable differences. On Quiz 1, we found
390 that participants' abilities to correctly answer questions about *Four Fundamental Forces* could be
391 predicted from their responses to questions about *Birth of Stars* ($OR = 1.896, \lambda_{LR} = 7.205, 95\% CI =$
392 $[6.224, 7.524], p = 0.039$) and similarly, that their ability to correctly answer *Birth of Stars*-related
393 questions could be predicted from their responses to *Four Fundamental Forces*-related questions
394 ($OR = 1.522, \lambda_{LR} = 6.448, 95\% CI = [5.656, 6.843], p = 0.043$). Given the results from our analyses
395 that included all questions and within-lecture predictions, we were surprised to find that the
396 knowledge estimates could reliably (if weakly) predict participants' performance across content
397 from different lectures. It is possible that this result reflects a combination of random guessing
398 prior to training (leading to a weak effect size), alongside some coarse-scale structured knowledge
399 that participants had about the content prior to watching either lecture. When we repeated
400 this analysis using questions from Quiz 2, we found participants' responses to *Four Fundamental*
401 *Forces*-related questions did *not* reliably predict their success on *Birth of Stars*-related questions
402 ($OR = 1.865, \lambda_{LR} = 3.205, 95\% CI = [3.027, 3.600], p = 0.125$), nor did their responses to *Birth of*
403 *Stars*-related questions reliably predict their success on *Four Fundamental Forces*-related questions
404 ($OR = 3.490, \lambda_{LR} = 3.266, 95\% CI = [3.033, 3.866], p = 0.094$). These "prediction failures" appear
405 to come from the fact that any signal derived from participants' knowledge about the content of
406 the *Birth of Stars* lecture (prior to watching it) is swamped by the much more dramatic increase
407 in their knowledge about the content of the *Four Fundamental Forces* (which they watched just
408 prior to taking Quiz 2). This is reflected in their Quiz 2 performance for questions about each
409 lecture (mean proportion correct for *Four Fundamental Forces*-related questions on Quiz 2: 0.77;
410 mean proportion correct for *Birth of Stars*-related questions on Quiz 2: 0.36). When we carried out

411 these across-lecture knowledge predictions using questions from Quiz 3 (when participants had
412 now viewed *both* lectures), we could again reliably predict success on questions about both *Four*
413 *Fundamental Forces* ($OR = 11.294$, $\lambda_{LR} = 11.055$, 95% CI = [9.126, 18.476], $p = 0.004$) and *Birth of*
414 *Stars* ($OR = 7.302$, $\lambda_{LR} = 7.068$, 95% CI = [6.490, 8.584], $p = 0.017$) using responses to questions
415 about the other lecture’s content. Across all three versions of these analyses, our results suggest
416 that (by and large) our knowledge estimates can reliably predict participants’ abilities to answer
417 individual quiz questions, distinguish between questions about similar content, and generalize
418 across content areas, provided that participants’ quiz responses reflect a minimum level of “real”
419 knowledge about both content on which these predictions are based and that for which they are
420 made. Our results also indicate some important limitations of our approach: if participants’ quiz
421 performance does not reflect what they know (e.g., when they “guess”), or if their knowledge is
422 not structured in a way that is reflected by the embedding space, then our knowledge estimates
423 will not be predictive of their performance.

424 That the knowledge predictions derived from the text embedding space reliably distinguish
425 between held-out correctly versus incorrectly answered questions (Fig. 6) suggests that spatial
426 relationships within this space can help explain what participants know. But how far does this
427 explanatory power extend? For example, suppose we know that a participant correctly answered a
428 question at embedding coordinate x . As we move farther away from x in the embedding space, how
429 does the likelihood that the participant knows about the content at a given location “fall off” with
430 distance? Conversely, suppose the participant instead answered that same question *incorrectly*.
431 Again, as we move farther away from x in the embedding space, how does the likelihood that the
432 participant does *not* know about a coordinate’s content change with distance? We reasoned that,
433 assuming our embedding space is capturing something about how individuals actually organize
434 their knowledge, a participant’s ability to answer questions embedded very close to x should
435 tend to be similar to their ability to answer the question embedded *at* x . Whereas at another
436 extreme, once we reach some sufficiently large distance from x , our ability to infer whether or
437 not a participant will correctly answer a question based on their ability to answer the question
438 at x should be no better than guessing based on their *overall* proportion of correctly answered



Figure 7: Knowledge falls off gradually in text embedding space. **A. Performance versus distance.** For each participant, for each correctly answered question (blue) or incorrectly answered question (red), we computed the proportion of correctly answered questions within a given distance of that question’s embedding coordinate. We used these proportions as a proxy for participants’ knowledge about the content within that region of the embedding space. We repeated this analysis for all questions and participants, and separately for each quiz (column). The black lines denote the average proportion correct across *all* questions included in the analysis at the given distance. **B. Maximum distance for which performance is reliably different from the average.** We used a bootstrap procedure (see *Estimating the “smoothness” of knowledge*) to estimate the point at which the blue and red lines in Panel A reliably diverged from the black line. We repeated this analysis separately for correctly and incorrectly answered questions from each quiz. **All panels.** Error ribbons denote bootstrap-estimated 95% confidence intervals.

439 questions. In other words, beyond the maximum distance at which the participant’s ability to
 440 answer the question at x is informative of their ability to answer a second question at location y ,
 441 then guessing the outcome at y based on x should be no more successful than guessing based on a
 442 measure that does not consider embedding space distance.

443 With these ideas in mind, we asked: conditioned on answering a question correctly, what
 444 proportion of all questions (within some radius, r , of that question’s embedding coordinate)
 445 were answered correctly? We plotted this proportion as a function of r . Similarly, we could
 446 ask, conditioned on answering a question incorrectly, how the proportion of correct responses
 447 changed with r . As shown in Figure 7, we found that quiz performance falls off smoothly with

448 distance, and the “rate” of the falloff does not appear to change across the different quizzes, as
449 measured by the distance at which performance becomes statistically indistinguishable from a
450 simple proportion correct score (see *Estimating the “smoothness” of knowledge*). This suggests that,
451 at least within the region of text embedding space covered by the questions our participants
452 answered (and as characterized using our topic model), the rate at which knowledge changes
453 with distance is relatively constant, even as participants’ overall level of knowledge varies across
454 quizzes or regions of the embedding space.

455 Knowledge estimates need not be limited to the content of the lectures. As illustrated in
456 Figure 8, our general approach to estimating knowledge from a small number of quiz questions
457 may be extended to *any* content, given its text embedding coordinate. To visualize how knowledge
458 “spreads” through text embedding space to content beyond the lectures participants watched, we
459 first fit a new topic model to the lectures’ sliding windows with $k = 100$ topics. Conceptually,
460 increasing the number of topics used by the model functions to increase the “resolution” of the
461 embedding space, providing a greater ability to estimate knowledge for content that is highly
462 similar to (but not precisely the same as) that contained in the two lectures. We note that we
463 used these 2D maps solely for visualization; all relevant comparisons, distance computations, and
464 statistical tests we report above were carried out in the original 15-dimensional space, using the
465 15-topic model. Aside from increasing the number of topics from 15 to 100, all other procedures
466 and model parameters were carried over from the preceding analyses. As in our other analyses,
467 we resampled each lecture’s topic trajectory to 1 Hz and projected each question into a shared text
468 embedding space.

469 We projected the resulting 100-dimensional topic vectors (for each second of video and each quiz
470 question) onto a shared 2-dimensional plane (see *Creating knowledge and learning map visualizations*).
471 Next, we sampled points from a 100×100 grid of coordinates that evenly tiled a rectangle enclos-
472 ing the 2D projections of the videos and questions. We used Equation 4 to estimate participants’
473 knowledge at each of these 10,000 sampled locations, and averaged these estimates across par-
474 ticipants to obtain an estimated average *knowledge map* (Fig. 8A). Intuitively, the knowledge map
475 constructed from a given quiz’s responses provides a visualization of how “much” participants



Figure 8: Mapping out the geometry of knowledge and learning. **A.** Average “knowledge maps” estimated using each quiz. Each map displays a 2D projection of the estimated knowledge about the content reflected by *all* regions of topic space (see *Creating knowledge and learning map visualizations*). The topic trajectories of the two lectures are indicated by dotted lines (blue: Lecture 1; green: Lecture 2), and the coordinates of each question are indicated by dots (light blue: Lecture 1-related; light green: Lecture 2-related; purple: general physics knowledge). Each map reflects an average across all participants. For individual participants’ maps, see Supplementary Figures 7, 8, and 9. **B.** Average “learning maps” estimated between each successive pair of quizzes. The learning maps follow the same general format as the knowledge maps in Panel A, but here the shading at each coordinate indicates the *difference* between the corresponding coordinates in the indicated pair of knowledge maps—i.e., how much the estimated knowledge “changed” between the two quizzes. Each map reflects an average across all participants. For individual participants’ maps, see Supplementary Figures 10 and 11. **C.** Word clouds for sampled points in topic space. Each word cloud displays the weighted blend of words underlying the topic proportions represented at the corresponding colored star’s location on the maps. In each word cloud, the words’ relative sizes correspond to their relative weights at the starred location, and their colors indicate their relative weights in *Four Fundamental Forces* (blue) versus *Birth of Stars* (green) lectures, on average, across all timepoints’ topic vectors.

476 knew about any content expressible by the fitted text embedding model at the point in time when
477 they completed that quiz.

478 Several features of the resulting knowledge maps are worth noting. The average knowledge
479 map estimated from Quiz 1 responses (Fig. 8A, leftmost map) shows that participants tended to
480 have relatively little knowledge about any parts of the text embedding space (i.e., the shading is
481 relatively dark everywhere). The knowledge map estimated from Quiz 2 responses shows a marked
482 increase in knowledge on the left side of the map (around roughly the same range of coordinates
483 traversed by the *Four Fundamental Forces* lecture, indicated by the dotted blue line). In other words,
484 participants' estimated increase in knowledge is localized to conceptual content that is nearby (i.e.,
485 related to) the content from the lecture they watched prior to taking Quiz 2. This localization is
486 non-trivial: these knowledge estimates are informed only by the embedded coordinates of the
487 *quiz questions*, not by the embeddings of either lecture (see Eqn. 4). Finally, the knowledge map
488 estimated from Quiz 3 responses shows a second increase in knowledge, localized to the region
489 surrounding the embedding of the *Birth of Stars* lecture participants watched immediately prior to
490 taking Quiz 3.

491 Another way of visualizing these content-specific increases in knowledge after participants
492 viewed each lecture is displayed in Figure 8B. Taking the point-by-point difference between the
493 knowledge maps estimated from responses to a successive pair of quizzes yields a *learning map*
494 that describes the *change* in knowledge estimates from one quiz to the next. These learning maps
495 highlight that the estimated knowledge increases we observed across maps were specific to the
496 regions around the embeddings of each lecture, in turn.

497 Because the 2D projection we used to construct the knowledge and learning maps is invertible,
498 we may gain additional insights into these maps' meanings by reconstructing the original high-
499 dimensional topic vector for any location on the map we are interested in. For example, this could
500 serve as a useful tool for an instructor looking to better understand which content areas a student
501 (or a group of students) knows well (or poorly). As a demonstration, we show the top-weighted
502 words from the blends of topics reconstructed from three example locations on the maps (Fig. 8C):
503 one point near the *Four Fundamental Forces* embedding (yellow), a second point near the *Birth of*

504 *Stars* embedding (orange), and a third point between the two lectures' embeddings (pink). As
505 shown in the word clouds in the panel, the top-weighted words at the example coordinate near the
506 *Four Fundamental Forces* embedding tended to be weighted more heavily by the topics expressed
507 in that lecture. Similarly, the top-weighted words at the example coordinate near the *Birth of Stars*
508 embedding tended to be weighted more heavily by the topics expressed in *that* lecture. And the
509 top-weighted words at the example coordinate between the two lectures' embeddings show a
510 roughly even mix of words most strongly associated with each lecture.

511 Discussion

512 We developed a computational framework that uses short multiple-choice quizzes to gain nuanced
513 insights into what learners know and how their knowledge changes with training. First, we show
514 that our approach can automatically match the conceptual knowledge probed by individual quiz
515 questions to the corresponding moments in lecture videos when those concepts were presented
516 (Fig. 4). Next, we demonstrate how we can estimate moment-by-moment “knowledge traces”
517 that reflect the degree of knowledge participants have about each video’s time-varying content,
518 and capture temporally specific increases in knowledge after viewing each lecture (Fig. 5). We
519 also show that these knowledge estimates can generalize to held-out questions (Fig. 6). Finally,
520 we use our framework to construct visual maps that provide snapshot estimates of how much
521 participants know about any concept within the scope of our text embedding model, and how
522 much their knowledge of those concepts changes with training (Fig. 8).

523 We view our work as making several contributions to the study of how people acquire con-
524 ceptual knowledge. First, from a methodological standpoint, our modeling framework provides
525 a systematic means of mapping out and characterizing knowledge in maps that have infinite (ar-
526 bitrarily many) numbers of coordinates, and of “filling out” those maps using relatively small
527 numbers of multiple choice quiz questions. Our experimental finding that we can use these maps
528 to predict responses to held-out questions has several psychological implications as well. For ex-
ample, concepts that are assigned to nearby coordinates by the text embedding model also appear

530 to be “known to a similar extent” (as reflected by participants’ responses to held-out questions;
531 Fig. 6). This suggests that participants also *conceptualize* similarly the content reflected by nearby
532 embedding coordinates. How participants’ knowledge falls off with spatial distance is captured
533 by the knowledge maps we infer from their quiz responses (e.g., Figs. 7, 8). In other words, our
534 study shows that knowledge about a given concept implies knowledge about related concepts,
535 and we also show how estimated knowledge falls off with distance in text embedding space.

536 In our study, we characterize the “coordinates” of participants’ knowledge using a relatively
537 simple “bag of words” text embedding model [LDA; 8]. More sophisticated text embedding mod-
538 els, such as transformer-based models [18, 55, 68, 71] can learn complex grammatical and semantic
539 relationships between words, higher-order syntactic structures, stylistic features, and more. We
540 considered using transformer-based models in our study, but we found that the text embeddings
541 derived from these models were surprisingly uninformative with respect to differentiating or oth-
542 erwise characterizing the conceptual content of the lectures and questions we used. We suspect
543 that this reflects a broader challenge in constructing models that are high-resolution within a given
544 domain (e.g., the domain of physics lectures and questions) *and* sufficiently broad so as to enable
545 them to cover a wide range of domains. For example, we found that the embeddings derived even
546 from much larger and more modern models like BERT [18], GPT [71], LLaMa [68], and others that
547 are trained on enormous text corpora, end up yielding poor resolution within the content space
548 spanned by individual course videos (Supp. Fig. 6). Whereas the LDA embeddings of the lectures
549 and questions are “near” each other (i.e., the convex hull enclosing the two lectures’ trajectories is
550 highly overlapping with the convex hull enclosing the questions’ embeddings), the BERT embed-
551 dings of the lectures and questions are instead largely distinct (top row of Supp. Fig. 6). The LDA
552 embeddings of the questions for each lecture and the corresponding lecture’s trajectory are also
553 similar. For example, as shown in Fig. 2C, the LDA embeddings for *Four Fundamental Forces* ques-
554 tions (blue dots) appear closer to the *Four Fundamental Forces* lecture trajectory (blue line), whereas
555 the LDA embeddings for *Birth of Stars* questions (green dots) appear closer to the *Birth of Stars*
556 lecture trajectory (green line). The BERT embeddings of the lectures and questions do not show
557 this property (Supp. Fig. 6). We also examined per-question “content matches” between individual

558 questions and individual moments of each lecture (Figs. 4, 6). The time series plot of individual
559 questions’ correlations are different from each other when computed using LDA (e.g., the traces
560 can be clearly visually separated), whereas the correlations computed from BERT embeddings of
561 different questions all look very similar. This tells us that LDA is capturing some differences in
562 content between the questions, whereas BERT is not. The time series plots of individual ques-
563 tions’ correlations have clear “peaks” when computed using LDA, but not when computed using
564 BERT. This tells us that LDA is capturing a “match” between the content of each question and a
565 relatively well-defined time window of the corresponding lectures. The BERT embeddings appear
566 to blur together the content of the questions versus specific moments of each lecture. Finally, we
567 also compared the pairwise correlations between embeddings of questions within versus across
568 content areas (i.e., content covered by the individual lectures, lecture-specific questions, and by the
569 “general physics knowledge” questions). The LDA embeddings show a strong contrast between
570 same-content embeddings versus across-content embeddings. In other words, the embeddings of
571 questions about the *Four Fundamental Forces* material are highly correlated with the embeddings of
572 the *Four Fundamental Forces* lecture, but not with the embeddings of *Birth of Stars*, questions about
573 *Birth of Stars*, or general physics knowledge questions. We see a similar pattern with the LDA
574 embeddings of the *Birth of Stars* questions (Fig. 3, Supp. Fig. 2). In contrast, the BERT embeddings
575 are all highly correlated with each other (Supp. Fig. 6). Taken together, these comparisons illus-
576 trate how LDA (trained on the specific content in question) provides both coverage of the requisite
577 material and specificity at the level of the content covered by individual questions. BERT, on the
578 other hand, essentially assigns both lectures and all of the questions (which are all broadly about
579 “physics”) into a tiny region of its embedding space, thereby blurring out meaningful distinctions
580 between different specific concepts covered by the lectures and questions. We note that these are
581 not criticisms of BERT (or other large language models trained on large and diverse corpora).
582 Rather, our point is that simple fine-tuned models trained on a relatively small but specialized
583 corpus can outperform much more complicated models trained on much larger corpora, when we
584 are specifically interested in capturing subtle conceptual differences at the level of a single course
585 lecture or question. Of course if our goal had been to find a model that generalized to many

586 different content areas, we would expect our approach to perform comparatively poorly relative to
587 BERT or other much larger models. We suggest that bridging the tradeoff between high resolution
588 within each content area versus the ability to generalize to many different content areas will be an
589 important challenge for future work in this domain.

590 Another application for large language models that does *not* require explicitly modeling the
591 content of individual lectures or questions is to leverage the models' abilities to generate text. For
592 example, generative text models like ChatGPT [55] and LLaMa [68] are already being used to build
593 a new generation of interactive tutoring systems [e.g., 44]. Unlike the approach we have taken here,
594 these generative text model-based systems do not explicitly model what learners know, or how
595 their knowledge changes over time with training. One could imagine building a hybrid system
596 that combines the best of both worlds: a large language model that can *generate* text, combined
597 with a smaller model that can *infer* what learners know and how their knowledge changes over
598 time. Such a hybrid system could potentially be used to build the next generation of interactive
599 tutoring systems that are able to adapt to learners' needs in real time, and that are able to provide
600 more nuanced feedback about what learners know and what they do not know.

601 At the opposite end of the spectrum from large language models, one could also imagine
602 *simplifying* some aspects of our LDA-based approach by computing simple word overlap metrics.
603 For example, the Jaccard similarity between text A and B is computed as the number of unique
604 words in the intersection of words from A and B divided by the number of unique words in the
605 union of words from A and B . In a supplementary analysis (Supp. Fig. 5), we compared the
606 LDA-based question-lecture matches we reported in Figure 4 with the Jaccard similarities between
607 each question and each sliding window of text from the corresponding lecture. As shown in
608 Supplementary Figure 5, this simple word-matching approach does not appear to capture the same
609 level of specificity as the LDA-based approach. Whereas the LDA-based approach often yields a
610 clear peak in the time series of correlations between each question and the corresponding lecture,
611 the Jaccard similarity-based approach does not. Furthermore, these LDA-based matches appear
612 to capture conceptual overlaps between the questions and lectures (Supp. Tab. 3), whereas simple
613 word matching does not. For example, one of the example questions examined in Supplementary

614 Figure 5 asks “Which of the following occurs as a cloud of atoms gets more dense?” The LDA-based
615 matches identify lecture timepoints where the relevant *topics* are discussed (e.g., when words like
616 “cloud,” “atom,” “dense,” etc., are mentioned *together*). The Jaccard similarity-based matches,
617 on the other hand, are strong when *any* of these words are mentioned, even if they do not occur
618 together.

619 We view our approach as occupying a sort of “sweet spot,” between much larger language
620 models and simple word matching-based approaches, that enables us to capture the relevant
621 conceptual content of course materials at an appropriate semantic scale. Our approach enables us
622 to accurately and consistently identify each question’s content in a way that also matches up with
623 what is presented in the lectures. In turn, this enables us to construct accurate predictions about
624 participants’ knowledge of the conceptual content tested by held-out questions (Fig. 6).

625 One limitation of our approach is that topic models contain no explicit internal representations
626 of more complex aspects of “knowledge,” like knowledge graphs, dependencies or associations
627 between concepts, causality, and so on. These representations might (in principle) be added
628 as extensions to our approach to more accurately and precisely capture, characterize, and track
629 learners’ knowledge. However, modeling these aspects of knowledge will likely require substantial
630 additional research effort.

631 Within the past several years, the global pandemic forced many educators to suddenly adapt to
632 teaching remotely [35, 51, 64, 72]. This change in world circumstances is happening alongside (and
633 perhaps accelerating) geometric growth in the availability of high-quality online courses from plat-
634 forms such as Khan Academy [36], Coursera [73], EdX [38], and others [60]. Continued expansion
635 of the global internet backbone and improvements in computing hardware have also facilitated
636 improvements in video streaming, enabling videos to be easily shared and viewed by increasingly
637 large segments of the world’s population. This exciting time for online course instruction provides
638 an opportunity to re-evaluate how we, as a global community, educate ourselves and each other.
639 For example, we can ask: what defines an effective course or training program? Which aspects of
640 teaching might be optimized and/or augmented by automated tools? How and why do learning
641 needs and goals vary across people? How might we lower barriers to receiving a high-quality

642 education?

643 Alongside these questions, there is a growing desire to extend existing theories beyond the
644 domain of lab testing rooms and into real classrooms [34]. In part, this has led to a recent
645 resurgence of “naturalistic” or “observational” experimental paradigms that attempt to better
646 reflect more ethologically valid phenomena that are more directly relevant to real-world situations
647 and behaviors [52]. In turn, this has brought new challenges in data analysis and interpretation. A
648 key step towards solving these challenges will be to build explicit models of real-world scenarios
649 and how people behave in them (e.g., models of how people learn conceptual content from real-
650 world courses, as in our current study). A second key step will be to understand which sorts
651 of signals derived from behaviors and/or other measurements [e.g., neurophysiological data; 4,
652 19, 49, 53, 57] might help to inform these models. A third major step will be to develop and
653 employ reliable ways of evaluating the complex models and data that are a hallmark of naturalistic
654 paradigms.

655 Beyond specifically predicting what people *know*, the fundamental ideas we develop here also
656 relate to the notion of “theory of mind” of other individuals [26, 32, 48]. Considering others’ unique
657 perspectives, prior experiences, knowledge, goals, etc., can help us to more effectively interact and
658 communicate [58, 63, 67]. One could imagine future extensions of our work (e.g., analogous to
659 the knowledge and learning maps shown in Fig. 8), that attempt to characterize how well-aligned
660 different people’s knowledge bases or backgrounds are. In turn, this might be used to model how
661 knowledge (or other forms of communicable information) flows not just between teachers and
662 students, but between friends having a conversation, individuals on a first date, participants at
663 a business meeting, doctors and patients, experts and non-experts, political allies or adversaries,
664 and more. For example, the extent to which two people’s knowledge maps “match” or “align” in
665 a given region of text embedding space might serve as a predictor of how effectively they will be
666 able to communicate about the corresponding conceptual content.

667 Ultimately, our work suggests a rich new line of questions about the geometric “form” of
668 knowledge, how knowledge changes over time, and how we might map out the full space of
669 what an individual knows. Our finding that detailed estimates about knowledge may be obtained

670 from short quizzes shows one way that traditional approaches to evaluation in education may be
671 extended. We hope that these advances might help pave the way for new approaches to teaching
672 or delivering educational content that are tailored to individual students' learning needs and goals.

673 Materials and methods

674 Participants

675 We enrolled a total of 50 Dartmouth undergraduate students in our study. Participants received
676 optional course credit for enrolling. We asked each participant to complete a demographic survey
677 that included questions about their age, gender, native spoken language, ethnicity, race, hearing,
678 color vision, sleep, coffee consumption, level of alertness, and several aspects of their educational
679 background and prior coursework.

680 Participants' ages ranged from 18 to 22 years (mean: 19.52 years; standard deviation: 1.09
681 years). A total of 15 participants reported their gender as male and 35 participants reported their
682 gender as female. A total of 49 participants reported their native language as "English" and 1
683 reported having another native language. A total of 47 participants reported their ethnicity as
684 "Not Hispanic or Latino" and three reported their ethnicity as "Hispanic or Latino." Participants
685 reported their races as White (32 participants), Asian (14 participants), Black or African American
686 (5 participants), American Indian or Alaska Native (1 participant), and Native Hawaiian or Other
687 Pacific Islander (1 participant). (Note that some participants selected multiple racial categories.)

688 A total of 49 participants reporting having normal hearing and 1 participant reported having
689 some hearing impairment. A total of 49 participants reported having normal color vision and 1
690 participant reported being color blind. Participants reported having had, on the night prior to
691 testing, 2–4 hours of sleep (1 participant), 4–6 hours of sleep (9 participants), 6–8 hours of sleep (35
692 participants), or 8+ hours of sleep (5 participants). They reported having consumed, on the same
693 day and leading up to their testing session, 0 cups of coffee (38 participants), 1 cup of coffee (10
694 participants), 3 cups of coffee (1 participant), or 4+ cups of coffee (1 participant).

695 No participants reported that their focus was currently impaired (e.g., by drugs or alcohol).
696 Participants reported their current level of alertness, and we converted their responses to numerical
697 scores as follows: "very sluggish" (-2), "a little sluggish" (-1), "neutral" (0), "fairly alert" (1), and
698 "very alert" (2). Across all participants, a range of alertness levels were reported (range: -2–1;
699 mean: -0.10; standard deviation: 0.84).

700 Participants reported their undergraduate major(s) as "social sciences" (28 participants), "nat-
701 ural sciences" (16 participants), "professional" (e.g., pre-med or pre-law; 8 participants), "mathe-
702 matics and engineering" (7 participants), "humanities" (4 participants), or "undecided" (3 partici-
703 pants). Note that some participants selected multiple categories for their undergraduate major(s).
704 We also asked participants about the courses they had taken. In total, 45 participants reported hav-
705 ing taken at least one Khan Academy course in the past, and 5 reported not having taken any Khan
706 Academy courses. Of those who reported having watched at least one Khan Academy course,
707 7 participants reported having watched 1–2 courses, 11 reported having watched 3–5 courses, 8
708 reported having watched 5–10 courses, and 19 reported having watched 10 or more courses. We
709 also asked participants about the specific courses they had watched, categorized under different
710 subject areas. In the "Mathematics" area, participants reported having watched videos on AP
711 Calculus AB (21 participants), Precalculus (17 participants), Algebra 2 (14 participants), AP Cal-
712 culus BC (12 participants), Trigonometry (11 participants), Algebra 1 (10 participants), Geometry
713 (8 participants), Pre-algebra (7 participants), Multivariable Calculus (5 participants), Differential
714 Equations (5 participants), Statistics and Probability (4 participants), AP Statistics (2 participants),
715 Linear Algebra (2 participants), Early Math (1 participant), Arithmetic (1 participant), and other
716 videos not listed in our survey (5 participants). In the "Science and engineering" area, participants
717 reported having watched videos on Chemistry, AP Chemistry, or Organic Chemistry (21 partici-
718 pants); Physics, AP Physics I, or AP Physics II (18 participants); Biology, AP Biology; or High
719 school Biology (15 participants); Health and Medicine (1 participant); or other videos not listed
720 in our survey (5 participants). We also asked participants whether they had specifically seen the
721 videos used in our experiment. Of the 45 participants who reported having having taken at least
722 one Khan Academy course in the past, 44 participants reported that they had not watched the *Four*

723 Fundamental Forces video, and 1 participant reported that they were not sure whether they had
724 watched it. All participants reported that they had not watched the *Birth of Stars* video. When
725 we asked participants about non-Khan Academy online courses, they reported having watched
726 or taken courses on Mathematics (15 participants), Science and engineering (11 participants), Test
727 preparation (9 participants), Economics and finance (3 participants), Arts and humanities (2 partic-
728 ipants), Computing (2 participants), and other categories not listed in our survey (17 participants).
729 Finally, we asked participants about in-person courses they had taken in different subject areas.
730 They reported taking courses in Mathematics (38 participants), Science and engineering (37 par-
731 ticipants), Arts and humanities (34 participants), Test preparation (27 participants), Economics
732 and finance (26 participants), Computing (14 participants), College and careers (7 participants), or
733 other courses not listed in our survey (6 participants).

734 **Experiment**

735 We hand-selected two course videos from the Khan Academy platform: *Four Fundamental Forces*
736 (an introduction to gravity, electromagnetism, the weak nuclear force, and the strong nuclear force;
737 duration: 10 minutes and 29 seconds) and *Birth of Stars* (an introduction to how stars are formed;
738 duration: 7 minutes and 57 seconds). All participants viewed the videos in the same order (i.e.,
739 *Four Fundamental Forces* followed by *Birth of Stars*).

740 We then hand-created 39 multiple-choice questions: 15 about the conceptual content of *Four*
741 *Fundamental Forces* (i.e., Lecture 1), 15 about the conceptual content of *Birth of Stars* (i.e., Lecture 2),
742 and 9 questions that tested for general conceptual knowledge about basic physics (covering material
743 that was not presented in either video). To help broaden the set of lecture-specific questions,
744 our team worked through each lecture in small segments to identify what each segment was
745 “about” conceptually, and then write a question about that concept. The general physics questions
746 were drawn our team’s prior coursework and areas of interest, along with internet searches and
747 brainstorming with the project team and other members of J.R.M.’s lab. Although we attempted to
748 design the questions to test “conceptual knowledge,” we note that estimating the specific “amount”
749 of conceptual understanding that each question “requires” to answer is somewhat subjective, and

750 might even come down to the “strategy” a given participant uses to answer the question at that
751 particular moment. The full set of questions and answer choices may be found in Supplementary
752 Table 1. The final set of questions (and response options) was reviewed and approved by J.R.M.
753 before we collected or analyzed the text or experimental data.

754 Over the course of the experiment, participants completed three 13-question multiple-choice
755 quizzes: the first before viewing Lecture 1, the second between Lectures 1 and 2, and the third
756 after viewing Lecture 2 (see Fig. 1). The questions appearing on each quiz, for each participant,
757 were randomly chosen from the full set of 39, with the constraints that (a) each quiz contained
758 exactly 5 questions about Lecture 1, 5 questions about Lecture 2, and 3 questions about general
759 physics knowledge, and (b) each question appear exactly once for each participant. The orders of
760 questions on each quiz, and the orders of answer options for each question, were also randomized.
761 We obtained informed consent from all participants, and our experimental protocol was approved
762 by the Committee for the Protection of Human Subjects at Dartmouth College. We used this
763 experiment to develop and test our computational framework for estimating knowledge and
764 learning.

765 **Analysis**

766 **Statistics**

767 All of the statistical tests performed in our study were two-sided. The 95% confidence intervals
768 we reported for each correlation were estimated by generating 10,000 bootstrap distributions of
769 correlation coefficients by sampling (with replacement) from the observed data.

770 **Constructing text embeddings of multiple lectures and questions**

771 We adapted an approach we developed in prior work [29] to embed each moment of the two
772 lectures and each question in our pool in a common representational space. Briefly, our approach
773 uses a topic model [Latent Dirichlet Allocation; 8] trained on a set of documents, to discover a set
774 of k “topics” or “themes.” Formally, each topic is defined as a distribution of weights over words

775 in the model’s vocabulary (i.e., the union of all unique words, across all documents, excluding
776 “stop words.”). Conceptually, each topic is intended to give larger weights to words that are
777 semantically related (as inferred from their tendency to co-occur in the same document). After
778 fitting a topic model, each document in the training set, or any *new* document that contains at
779 least some of the words in the model’s vocabulary, may be represented as a k -dimensional vector
780 describing how much the document (most probably) reflects each topic. To select an appropriate k
781 for our model, as a starting point, we identified the minimum number of topics that yielded at least
782 one “unused” topic (i.e., in which all words in the vocabulary were assigned uniform weights)
783 after training. This indicated that the number of topics was sufficient to capture the set of latent
784 themes present in the two lectures (from which we constructed our document corpus, as described
785 below). We found this value to be $k = 15$ topics. We found that with a limited number of additional
786 adjustments following Boyd-Graber et al. [9], such as removing corpus-specific stop-words, the
787 model yielded (subjectively) sensible and coherent topics. The distribution of weights over words
788 in the vocabulary for each discovered topic is shown in Supplementary Figure 1, and each topic’s
789 top-weighted words may be found in Supplementary Table 2.

790 As illustrated in Figure 2A, we start by building up a corpus of documents using overlapping
791 sliding windows that span each video’s transcript. Khan Academy provides professionally created,
792 manual transcriptions of all videos for closed captioning. However, such transcripts would not
793 be readily available in all contexts to which our framework could potentially be applied. Khan
794 Academy videos are hosted on the YouTube platform, which additionally provides automated
795 captions. We opted to use these automated transcripts [which, in prior work, we have found to be
796 of sufficiently near-human quality to yield reliable data in behavioral studies; 74] when developing
797 our framework in order to make it more directly extensible and adaptable by others in the future.

798 We fetched these automated transcripts using the `youtube-transcript-api` Python pack-
799 age [17]. The transcripts consisted of one timestamped line of text for every few seconds (mean:
800 2.34 s; standard deviation: 0.83 s) of spoken content in the video (i.e., corresponding to each indi-
801 vidual caption that would appear on-screen if viewing the lecture via YouTube, and when those
802 lines would appear). We defined a sliding window length of (up to) $w = 30$ transcript lines, and

803 assigned each window a timestamp corresponding to the midpoint between the timestamps for its
804 first and last lines. This w parameter was chosen to match the same number of words per sliding
805 window (rounded to the nearest whole word, and before preprocessing) as the sliding windows
806 we defined in our prior work [29; i.e., 185 words per sliding window].

807 These sliding windows ramped up and down in length at the beginning and end of each
808 transcript, respectively. In other words, each transcript’s first sliding window covered only its first
809 line, the second sliding window covered the first two lines, and so on. This ensured that each line
810 from the transcripts appeared in the same number (w) of sliding windows. We next performed a
811 series of standard text preprocessing steps: normalizing case, lemmatizing, removing punctuation
812 and removing stop-words. We constructed our corpus of stop words by augmenting the Natural
813 Language Toolkit [NLTK; 5] English stop word list with the following additional words, selected
814 using one of the approaches suggested by Boyd-Graber et al. [9]: “actual,” “actually,” “also,” “bit,”
815 “could,” “e,” “even,” “first,” “follow,” “following,” “four,” “let,” “like,” “mc,” “really,” “saw,”
816 “see,” “seen,” “thing,” and “two.” This yielded sliding windows with an average of 73.8 remaining
817 words, and lasting for an average of 62.22 seconds. We treated the text from each sliding window
818 as a single “document,” and combined these documents across the two videos’ windows to create
819 a single training corpus for the topic model.

820 After fitting a topic model to the two videos’ transcripts, we could use the trained model to
821 transform arbitrary (potentially new) documents into k -dimensional topic vectors. A convenient
822 property of these topic vectors is that documents that reflect similar blends of topics (i.e., documents
823 that reflect similar themes, according to the model) will yield similar coordinates (in terms of
824 correlation, cosine similarity, Kullback-Leibler divergence, Euclidean distance, or other geometric
825 measures). In general, the similarity between different documents’ topic vectors may be used to
826 characterize the similarity in conceptual content between the documents.

827 We transformed each sliding window’s text into a topic vector, and then used linear interpolation
828 (independently for each topic dimension) to resample the resulting time series to one vector
829 per second. We also used the fitted model to obtain topic vectors for each question in our pool (see
830 Supp. Tab. 1). Taken together, we obtained a *trajectory* for each video, describing its path through

831 topic space, and a single coordinate for each question (Fig. 2C). Embedding both videos and all of
 832 the questions using a common model enables us to compare the content from different moments
 833 of videos, compare the content across videos, and estimate potential associations between specific
 834 questions and specific moments of video.

835 **Estimating dynamic knowledge traces**

836 We used the following equation to estimate each participant’s knowledge about timepoint t of a
 837 given lecture, $\hat{k}(t)$:

$$\hat{k}(f(t, L)) = \frac{\sum_{i \in \text{correct}} \text{ncorr}(f(t, L), f(i, Q))}{\sum_{j=1}^N \text{ncorr}(f(t, L), f(j, Q))}, \quad (1)$$

838 where

$$\text{ncorr}(x, y) = \frac{\text{corr}(x, y) - \text{mincorr}}{\text{maxcorr} - \text{mincorr}}, \quad (2)$$

839 and where mincorr and maxcorr are the minimum and maximum correlations between any lecture
 840 timepoint and question, taken over all timepoints in the given lecture, and all five questions *about*
 841 that lecture appearing on the given quiz. We also define $f(s, \Omega)$ as the s^{th} topic vector from the set
 842 of topic vectors Ω . Here t indexes the set of lecture topic vectors, L , and i and j index the topic
 843 vectors of questions used to estimate the knowledge trace, Q . Note that “correct” denotes the set
 844 of indices of the questions the participant answered correctly on the given quiz.

845 Intuitively, $\text{ncorr}(x, y)$ is the correlation between two topic vectors (e.g., the topic vector from one
 846 timepoint in a lecture, x , and the topic vector for one question, y), normalized by the minimum and
 847 maximum correlations (across all timepoints t and questions Q) to range between 0 and 1, inclusive.
 848 Equation 1 then computes the weighted average proportion of correctly answered questions about
 849 the content presented at timepoint t , where the weights are given by the normalized correlations
 850 between timepoint t ’s topic vector and the topic vectors for each question. The normalization step
 851 (i.e., using ncorr instead of the raw correlations) ensures that every question contributes some
 852 non-negative amount to the knowledge estimate.

853 **Generalized linear mixed models**

854 In the set of analyses reported in Figure 6, we assessed whether estimates of participants' knowledge
855 at the embedding coordinates of individual quiz questions could be used to reliably predict their
856 ability to correctly answer those questions. In essence, we treated each question a given participant
857 answered on a given quiz as a "lecture" consisting of a single timepoint, and used Equation 1 to
858 estimate the participant's knowledge for its embedding coordinate based on their performance on
859 all *other* questions they answered on that same quiz ("All questions"; Fig. 6, top row). Additionally,
860 for each lecture-related question (i.e., excluding questions about general physics knowledge), we
861 computed analogous knowledge estimates based on all other questions the participant answered
862 on the same quiz about (1) the same lecture as the target question ("Within-lecture"; Fig. 6, middle
863 rows), and (2) the other of the two lectures ("Across-lecture"; Fig. 6, bottom rows).

864 In each version of this analysis (i.e., row in Fig. 6), and separately for each of the three quizzes
865 (i.e., column in Fig. 6), we then fit a generalized linear mixed model (GLMM) with a logistic link
866 function to the set of knowledge estimates for all questions that participants answered on the
867 given quiz. We implemented these models in R using the `lme4` package [3] and fit them following
868 guidance from Bates et al. [2] and Matuschek et al. [45]. Specifically, we initially fit each model
869 with the maximal random effects structure afforded by our design, which we identified as:

$$\text{accuracy} \sim \text{knowledge} + (\text{knowledge} | \text{participant}) + (\text{knowledge} | \text{question})$$

870 where "accuracy" is a binary value indicating whether each target question was answered cor-
871 rectly or incorrectly, "knowledge" is estimated knowledge at each target question's embedding
872 coordinate, "participant" is a unique identifier assigned to each participant, and "question" is a
873 unique identifier assigned to each quiz question. For models we fit using knowledge estimates for
874 target questions about multiple content areas (i.e., in the "All questions" version of the analysis),
875 we also included an additional random effect term, $(\text{knowledge} | \text{lecture})$, where "lecture" is a
876 categorical value denoting whether the target question was about *Four Fundamental Forces*, *Birth*
877 *of Stars*, or general physics knowledge. Note that with our coding scheme, identifiers for each

878 question are implicitly nested within levels of lecture and do not require explicit nesting in
879 our model formula. We then iteratively removed random effects from the maximal model until
880 it successfully converged with a full rank (i.e., non-singular) random effects variance-covariance
881 matrix.

882 To assess the predictive value of our knowledge estimates, we compared each GLMM’s ability
883 to discriminate between correctly and incorrectly answered questions to that of an analogous model
884 that did *not* consider estimated knowledge. Specifically, we used the same sets of observations
885 with which we fit each “full” model to fit a second “null” model, with the formula:

$$\text{accuracy} \sim (1 | \text{participant}) + (1 | \text{question})$$

886 where “accuracy”, “participant”, and “question” are as defined above. As with our full models,
887 the null models we fit for the “All questions” version of the analysis for each quiz contained an
888 additional term, $(1 | \text{lecture})$, where “lecture” is as defined above. We then compared each
889 full model to its reduced (null) equivalent using a likelihood-ratio test (LRT). Because the typical
890 asymptotic χ^2_d approximation of the null distribution for the LRT statistic (λ_{LR}) is anti-conservative
891 for models that differ in their random slope terms [25, 61, 66], we computed *p*-values for these
892 tests using a parametric bootstrapping procedure [27]. For each of 1,000 bootstraps, we used the
893 fitted null model to simulate a sample of observations of equal size to our original sample. We
894 then re-fit both the null and full models to this simulated sample and compared them via an LRT.
895 This yielded a distribution of λ_{LR} statistics we may expect to observe under our null hypothesis.
896 Following [14, 54], we computed a corrected *p*-value for our observed λ_{LR} as $\frac{r+1}{n+1}$, where r is the
897 number of simulated model comparisons that yielded a λ_{LR} greater than or equal to our observed
898 value and n is the number of simulations we ran (1,000).

899 **Estimating the “smoothness” of knowledge**

900 In the analysis reported in Figure 7A, we show how participants’ ability to correctly answer
901 quiz questions changes as a function of distance from a given correctly or incorrectly answered

902 reference question. We used a bootstrap-based approach to estimate the maximum distances over
903 which these proportions of correctly answered questions could be reliably distinguished from
904 participants' overall average proportion of correctly answered questions.

905 For each of 10,000 iterations, we drew a random subsample (with replacement) of 50 partic-
906 ipants from our dataset. Within each iteration, we first computed the 95% confidence interval
907 (CI) of the across-subsample-participants mean proportion correct on each of the three quizzes,
908 separately. To compute this interval for each quiz, we repeatedly (1,000 times) subsampled par-
909 ticipants (with replacement, from the outer subsample for the current iteration) and computed
910 the mean proportion correct of each of these inner subsamples. We then identified the 2.5th and
911 97.5th percentiles of the resulting distributions of 1,000 means. These three intervals (one for each
912 quiz) served as our thresholds for confidence that the proportion correct within a given distance
913 from a reference question was reliably different (at the $p < 0.05$ significance level) from the average
914 proportion correct across all questions on the given quiz.

915 Next, for each participant in the current subsample, and for each of the three quizzes they
916 completed (separately), we iteratively treated each of the 15 questions appearing on the given
917 quiz as the "reference" question. We constructed a series of concentric 15-dimensional "spheres"
918 centered on the reference question's embedding space coordinate, where each successive sphere's
919 radius increased by 0.01 (correlation distance) between 0 and 2, inclusive (i.e., tiling the range
920 of possible correlation distances with 201 spheres in total). We then computed the proportion
921 of questions enclosed within each sphere that the participant answered correctly, and averaged
922 these per-radius proportion correct scores across reference questions that were answered correctly,
923 and those that were answered incorrectly. This resulted in two number-of-spheres sequences of
924 proportion-correct scores for each subsample participant and quiz: one derived from correctly
925 answered reference questions, and one derived from incorrectly answered reference questions.

926 We computed the across-subsample-participants mean proportion correct for each radius value
927 (i.e., sphere) and "correctness" of reference question. This yielded two sequences of proportion-
928 correct scores for each quiz, analogous to the blue and red lines displayed in Figure 7A, but for
929 the present subsample. For each quiz, we then found the minimum distance from the reference

930 question (i.e., sphere radius) at which each of these two sequences of per-radius proportion correct
931 scores intersected the 95% confidence interval for the overall proportion correct (i.e., analogous to
932 the black error bands in Fig. 7A).

933 This resulted in two “intersection” distances for each quiz (for correctly answered and incor-
934 rectly answered reference questions). Repeating this full process for each of the 10,000 bootstrap
935 iterations output two distributions of intersection distances for each of the three quizzes. The
936 means and 95% confidence intervals for these distributions are plotted in Figure 7B.

937 **Creating knowledge and learning map visualizations**

938 An important feature of our approach is that, given a trained text embedding model and partic-
939 ipants’ quiz performance on each question, we can estimate their knowledge about *any* content
940 expressible by the embedding model—not solely the content explicitly probed by the quiz ques-
941 tions, or even appearing in the lectures. To visualize these estimates (Fig. 8, Supp. Figs. 7, 8, 9, 10,
942 and 11), we used Uniform Manifold Approximation and Projection [UMAP; 46, 47] to construct a
943 2D projection of the text embedding space. Whereas our main analyses used a 15-topic embedding
944 space, we used a 100-topic embedding space for these visualizations. This change in the number
945 of topics overcame an undesirable behavior in the UMAP embedding procedure, whereby embed-
946 ding coordinates for the 15-topic model tended to be “clumped” into separated clusters, rather
947 than forming a smooth trajectory through the 2D space. When we increased the number of topics
948 to 100, the embedding coordinates in the 2D space formed a smooth trajectory through the space,
949 with substantially less clumping (Fig. 8). Creating a “map” by sampling this 100-dimensional
950 space at high resolution to obtain an adequate set of topic vectors spanning the embedding space
951 would be computationally intractable. However, sampling a 2D grid is trivial.

952 At a high level, the UMAP algorithm obtains low-dimensional embeddings by minimizing
953 the cross-entropy between the pairwise (clustered) distances between the observations in their
954 original (e.g., 100-dimensional) space and the pairwise (clustered) distances in the low-dimensional
955 embedding space (in our approach, the embedding space is 2D). In our implementation, pairwise
956 distances in the original high-dimensional space were defined as 1 minus the correlation between

957 each pair of coordinates, and pairwise distances in the low-dimensional embedding space were
958 defined as the Euclidean distance between each pair of coordinates.

959 In our application, all of the coordinates we embedded were topic vectors, whose elements
960 are always non-negative and sum to one. Although UMAP is an invertible transformation at
961 the embedding locations of the original data, other locations in the embedding space will not
962 necessarily follow the same implicit “rules” as the original high-dimensional data. For example,
963 inverting an arbitrary coordinate in the embedding space might result in negative-valued vectors,
964 which are incompatible with the topic modeling framework. To protect against this issue, we
965 log-transformed the topic vectors prior to embedding them in the 2D space. When we inverted
966 the embedded vectors (e.g., to estimate topic vectors or word clouds, as in Fig. 8C), we passed
967 the inverted (log-transformed) values through the exponential function to obtain a vector of non-
968 negative values, and normalized them to sum to one.

969 After embedding both lectures’ topic trajectories and the topic vectors of every question, we
970 defined a rectangle enclosing the 2D projections of the lectures’ and quizzes’ embeddings. We then
971 sampled points from a regular 100×100 grid of coordinates that evenly tiled this enclosing rectangle.
972 We sought to estimate participants’ knowledge (and learning, i.e., changes in knowledge) at each
973 of the resulting 10,000 coordinates.

974 To generate our estimates, we placed a set of 39 radial basis functions (RBFs) throughout the
975 embedding space, centered on the 2D projections for each question (i.e., we included one RBF for
976 each question). At coordinate x , the value of an RBF centered on a question’s coordinate μ , is given
977 by:

$$\text{RBF}(x, \mu, \lambda) = \exp \left\{ -\frac{\|x - \mu\|^2}{\lambda} \right\}. \quad (3)$$

978 The λ term in the RBF equation controls the “smoothness” of the function, where larger values
979 of λ result in smoother maps. In our implementation we used $\lambda = 50$. Next, we estimated the
980 “knowledge” at each coordinate, x , using:

$$\hat{k}(x) = \frac{\sum_{i \in \text{correct}} \text{RBF}(x, q_i, \lambda)}{\sum_{j=1}^N \text{RBF}(x, q_j, \lambda)}. \quad (4)$$

981 Intuitively, Equation 4 computes the weighted proportion of correctly answered questions, where
982 the weights are given by how nearby (in the 2D space) each question is to the x . We also defined
983 *learning maps* as the coordinate-by-coordinate differences between any pair of knowledge maps.
984 Intuitively, learning maps reflect the *change* in knowledge across two maps.

985 **Author contributions**

986 Conceptualization: P.C.F., A.C.H., and J.R.M. Methodology: P.C.F., A.C.H., and J.R.M. Software:
987 P.C.F. Validation: P.C.F. Formal analysis: P.C.F. Resources: P.C.F., A.C.H., and J.R.M. Data curation:
988 P.C.F. Writing (original draft): J.R.M. Writing (review and editing): P.C.F., A.C.H., and J.R.M. Visu-
989 alization: P.C.F. and J.R.M. Supervision: J.R.M. Project administration: P.C.F. Funding acquisition:
990 J.R.M.

991 **Data availability**

992 All of the data analyzed in this manuscript may be found at <https://github.com/ContextLab/efficient-learning-khan>.
993

994 **Code availability**

995 All of the code for running our experiment and carrying out the analyses may be found at
996 <https://github.com/ContextLab/efficient-learning-khan>.

997 **Acknowledgements**

998 We acknowledge useful discussions, assistance in setting up an earlier (unpublished) version of
999 this study, and assistance with data collection efforts from Will Baxley, Max Bluestone, Daniel
1000 Carstensen, Kunal Jha, Caroline Lee, Lucy Owen, Xinming Xu, and Kirsten Ziman. Our work was

1001 supported in part by NSF CAREER Award Number 2145172 to J.R.M. The content is solely the
1002 responsibility of the authors and does not necessarily represent the official views of our supporting
1003 organizations. The funders had no role in study design, data collection and analysis, decision to
1004 publish, or preparation of the manuscript.

1005 **References**

- 1006 [1] Ashby, F. G. and Maddox, W. T. (2005). Human category learning. *Annual Review of Psychology*,
1007 56:149–178.
- 1008 [2] Bates, D., Kliegl, R., Vasishth, S., and Baayen, H. (2015a). Parsimonious mixed models. *arXiv*,
1009 1506.04967.
- 1010 [3] Bates, D., Mächler, M., Bolker, B., and Walker, S. (2015b). Fitting linear mixed-effects models
1011 using lme4. *Journal of Statistical Software*, 67(1):1–48.
- 1012 [4] Bevilacqua, D., Davidesco, I., Wan, L., and Chaloner, K. (2019). Brain-to-brain synchrony and
1013 learning outcomes vary by student-teacher dynamics: evidence from a real-world classroom
1014 electroencephalography study. *Journal of Cognitive Neuroscience*, 31(3):401–411.
- 1015 [5] Bird, S., Klein, E., and Loper, E. (2009). *Nature language processing with Python: analyzing text
1016 with the natural language toolkit*. Reilly Media, Inc.
- 1017 [6] Blaye, A., Bernard-Peyron, V., Paour, J.-L., and Bonthoux, F. (2006). Category flexibility in chil-
1018 dren: distinguishing response flexibility from conceptual flexibility; the protracted development
1019 of taxonomic representations. *European Journal of Developmental Psychology*, 3(2):163–188.
- 1020 [7] Blei, D. M. and Lafferty, J. D. (2006). Dynamic topic models. In *Proceedings of the International
1021 Conference on Machine Learning*, pages 113–120, New York, NY. Association for Computing
1022 Machinery.
- 1023 [8] Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine
1024 Learning Research*, 3:993–1022.

- 1025 [9] Boyd-Graber, J., Mimno, D., and Newman, D. (2014). Care and feeding of topic models:
1026 problems, diagnostics, and improvements. In Airoldi, E. M., Blei, D. M., Erosheva, E. A., and
1027 Fienberg, S. E., editors, *Handbook of Mixed Membership Models and Their Applications*. CRC Press.
- 1028 [10] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A.,
1029 Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child,
1030 R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M.,
1031 Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei,
1032 D. (2020). Language models are few-shot learners. *arXiv*, 2005.14165.
- 1033 [11] Caramazza, A. and Mahon, B. Z. (2003). The organization of conceptual knowledge: the
1034 evidence from category-specific semantic deficits. *Trends in Cognitive Sciences*, 7(8):354–361.
- 1035 [12] Cer, D., Yang, Y., Kong, S. Y., Hua, N., Limtiaco, N., John, R. S., Constant, N., Guajardo-
1036 Cespedes, M., Yuan, S., Tar, C., Sung, Y.-H., Strope, B., and Kurzweil, R. (2018). Universal
1037 sentence encoder. *arXiv*, 1803.11175.
- 1038 [13] Constantinescu, A. O., O'Reilly, J. X., and Behrens, T. E. J. (2016). Organizing conceptual
1039 knowledge in humans with a gridlike code. *Science*, 352(6292):1464–1468.
- 1040 [14] Davison, A. C. and Hinkley, D. V. (1997). *Bootstrap Methods and their Application*. Cambridge
1041 Series in Statistical and Probabilistic Mathematics. Cambridge University Press.
- 1042 [15] Deacon, D., Grose-Fifer, J., Yang, C. M., Stanick, V., Hewitt, S., and Dynowska, A. (2004).
1043 Evidence for a new conceptualization of semantic representation in the left and right cerebral
1044 hemispheres. *Cortex*, 40(3):467–478.
- 1045 [16] Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., and Harshman, R. (1990).
1046 Indexing by latent semantic analysis. *Journal of the American Society for Information Science*,
1047 41(6):391–407.
- 1048 [17] Depoix, J. (2018). YouTube transcript API. <https://github.com/jdepoix/youtube-transcript-api>.

- 1050 [18] Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). BERT: pre-training of deep
1051 bidirectional transformers for language understanding. *arXiv*, 1810.04805.
- 1052 [19] Dikker, S., Wan, L., Davidesco, I., Kaggen, L., Oostrik, M., McClintock, J., Rowland, J.,
1053 Michalareas, G., van Bavel, J. J., Ding, M., and Poeppel, D. (2017). Brain-to-brain synchrony
1054 tracks real-world dynamic group interactions in the classroom. *Current Biology*, 27(9):1375–1380.
- 1055 [20] Estes, W. K. (1986a). Array models for category learning. *Cognitive Psychology*, 18(4):500–549.
- 1056 [21] Estes, W. K. (1986b). Memory storage and retrieval processes in category learning. *Journal of*
1057 *Experimental Psychology: General*, 115:155–174.
- 1058 [22] Fisher, R. A. (1922). On the mathematical foundations of theoretical statistics. *Philosophical*
1059 *Transactions of the Royal Society A*, 222(602):309–368.
- 1060 [23] Gallagher, J. J. (2000). Teaching for understanding and application of science knowledge.
1061 *School Science and Mathematics*, 100(6):310–318.
- 1062 [24] Gluck, M. A., Shohamy, D., and Myers, C. E. (2002). How do people solve the “weather
1063 prediction” task? individual variability in strategies for probabilistic category learning. *Learning*
1064 *and Memory*, 9:408–418.
- 1065 [25] Goldman, N. and Whelan, S. (2000). Statistical Tests of Gamma-Distributed Rate Heterogeneity
1066 in Models of Sequence Evolution in Phylogenetics. *Molecular Biology and Evolution*, 17(6):975–978.
- 1067 [26] Goldstein, T. R. and Winner, E. (2012). Enhancing empathy and theory of mind. *Journal of*
1068 *Cognition and Development*, 13(1):19–37.
- 1069 [27] Halekoh, U. and Højsgaard, S. (2014). A Kenward-Roger Approximation and Parametric
1070 Bootstrap Methods for Tests in Linear Mixed Models – The R Package pbkrtest. *Journal of*
1071 *Statistical Software*, 59(9):1–32.
- 1072 [28] Hall, R. and Greeno, J. (2008). *21st century education: A reference handbook*, chapter Conceptual
1073 learning, pages 212–221. Sage Publications.

- 1074 [29] Heusser, A. C., Fitzpatrick, P. C., and Manning, J. R. (2021). Geometric models reveal behav-
1075 ioral and neural signatures of transforming experiences into memories. *Nature Human Behaviour*,
1076 5:905–919.
- 1077 [30] Huebner, P. A. and Willits, J. A. (2018). Structured semantic knowledge can emerge au-
1078 tomatically from predicting word sequences in child-directed speech. *Frontiers in Psychology*,
1079 9:doi.org/10.3389/fpsyg.2018.00133.
- 1080 [31] Hulbert, J. C. and Norman, K. A. (2015). Neural differentiation tracks improved recall of com-
1081 peting memories following interleaved study and retrieval practice. *Cerebral Cortex*, 25(10):3994–
1082 4008.
- 1083 [32] Kanske, P., Böckler, A., and Singer, T. (2015). Models, mechanisms and moderators dissociating
1084 empathy and theory of mind. In *Social Behavior From Rodents to Humans*, pages 193–206. Springer.
- 1085 [33] Katona, G. (1940). *Organizing and memorizing: studies in the psychology of learning and teaching*.
1086 Columbia University Press.
- 1087 [34] Kaufman, D. M. (2003). Applying educational theory in practice. *British Medical Journal*,
1088 326(7382):213–216.
- 1089 [35] Kawasaki, H., Yamasaki, S., Masuoka, Y., Iwasa, M., Fukita, S., and Matsuyama, R. (2021).
1090 Remote teaching due to COVID-19: an exploration of its effectiveness and issues. *International
1091 Journal of Environmental Research and Public Health*, 18(5):2672.
- 1092 [36] Khan, S. (2004). *The Khan Academy*. Salman Khan.
- 1093 [37] Kintsch (1970). *Learning, memory, and conceptual processes*. Wiley.
- 1094 [38] Kolowich, S. (2013). How EdX plans to earn, and share, revenue from its free online courses.
1095 *The Chronicle of Higher Education*, 21:1–5.
- 1096 [39] Landauer, T. K. and Dumais, S. T. (1997). A solution to Plato’s problem: the latent semantic
1097 analysis theory of acquisition, induction, and representation of knowledge. *Psychological Review*,
1098 104:211–240.

- 1099 [40] Lee, H. and Chen, J. (2022). Predicting memory from the network structure of naturalistic
1100 events. *Nature Communications*, 13(4235):doi.org/10.1038/s41467-022-31965-2.
- 1101 [41] Macellan, E. (2005). Conceptual learning: the priority for higher education. *British Journal of
1102 Educational Studies*, 53(2):129–147.
- 1103 [42] Manning, J. R. (2021). Episodic memory: mental time travel or a quantum “memory wave”
1104 function? *Psychological Review*, 128(4):711–725.
- 1105 [43] Manning, J. R. (2023). Context reinstatement. In Kahana, M. J. and Wagner, A. D., editors,
1106 *Handbook of Human Memory*. Oxford University Press.
- 1107 [44] Manning, J. R., Menjunatha, H., and Kording, K. (2023). Chatify: A Jupyter extension
1108 for adding LLM-driven chatbots to interactive notebooks. [https://github.com/ContextLab/
1109 chatify](https://github.com/ContextLab/chatify).
- 1110 [45] Matuschek, H., Kliegl, R., Vasishth, S., Baayen, H., and Bates, D. (2017). Balancing type i error
1111 and power in linear mixed models. *Journal of Memory and Language*, 94:305–315.
- 1112 [46] McInnes, L., Healy, J., and Melville, J. (2018a). UMAP: Uniform manifold approximation and
1113 projection for dimension reduction. *arXiv*, 1802(03426).
- 1114 [47] McInnes, L., Healy, J., Saul, N., and Großberger, L. (2018b). UMAP: Uniform Manifold
1115 Approximation and Projection. *Journal of Open Source Software*, 3(29):861.
- 1116 [48] Meltzoff, A. N. (2011). Social cognition and the origins of imitation, empathy, and theory of
1117 mind. In *The Wiley-Blackwell Handbook of Childhood Cognitive Development*. Wiley-Blackwell.
- 1118 [49] Meshulam, M., Hasenfratz, L., Hillman, H., Liu, Y. F., Nguyen, M., Norman, K. A., and Hasson,
1119 U. (2020). Neural alignment predicts learning outcomes in students taking an introduction to
1120 computer science course. *Nature Communications*, 12(1922):doi.org/10.1038/s41467-021-22202-3.
- 1121 [50] Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word represen-
1122 tations in vector space. *arXiv*, 1301.3781.

- 1123 [51] Moser, K. M., Wei, T., and Brenner, D. (2021). Remote teaching during COVID-19: implications
1124 from a national survey of language educators. *System*, 97:102431.
- 1125 [52] Nastase, S. A., Goldstein, A., and Hasson, U. (2020). Keep it real: rethinking the primacy of
1126 experimental control in cognitive neuroscience. *NeuroImage*, 15(222):117254–117261.
- 1127 [53] Nguyen, M., Chang, A., Micciche, E., Meshulam, M., Nastase, S. A., and Hasson, U. (2022).
1128 Teacher-student neural coupling during teaching and learning. *Social Cognitive and Affective*
1129 *Neuroscience*, 17(4):367–376.
- 1130 [54] North, B. V., Curtis, D., and Sham, P. C. (2002). A note on the calculation of empirical p values
1131 from monte carlo procedures. *American Journal of Human Genetics*, 71(2):439–441.
- 1132 [55] OpenAI (2023). ChatGPT. <https://chat.openai.com>.
- 1133 [56] Piantadosi, S. T. and Hill, F. (2022). Meaning without reference in large language models.
1134 *arXiv*, 2208.02957.
- 1135 [57] Poulsen, A. T., Kamronn, S., Dmochowski, J., Parra, L. C., and Hansen, L. K. (2017). EEG
1136 in the classroom: synchronised neural recordings during video presentation. *Scientific Reports*,
1137 7:43916.
- 1138 [58] Ratka, A. (2018). Empathy and the development of affective skills. *American Journal of*
1139 *Pharmaceutical Education*, 82(10):doi.org/10.5688/ajpe7192.
- 1140 [59] Reilly, D. L., Cooper, L. N., and Elbaum, C. (1982). A neural model for category learning.
1141 *Biological Cybernetics*, 45(1):35–41.
- 1142 [60] Rhoads, R. A., Berdan, J., and Toven-Lindsey, B. (2013). The open courseware movement in
1143 higher education: unmasking power and raising questions about the movement’s democratic
1144 potential. *Educational Theory*, 63(1):87–110.
- 1145 [61] Scheipl, F., Greven, S., and Küchenhoff, H. (2008). Size and power of tests for a zero random
1146 effect variance or polynomial regression in additive and linear mixed models. *Computational*
1147 *Statistics & Data Analysis*, 52(7):3283–3299.

- 1148 [62] Scott, P., Asoko, H., and Leach, J. (2007). *Handbook of research on science education*, chapter
1149 Student conceptions and conceptual learning in science. Routledge.
- 1150 [63] Shao, Y. N., Sun, H. M., Huang, J. W., Li, M. L., Huang, R. R., and Li, N. (2018). Simulation-
1151 based empathy training improves the communication skills of neonatal nurses. *Clinical Simula-*
1152 *tion in Nursing*, 22:32–42.
- 1153 [64] Shim, T. E. and Lee, S. Y. (2020). College students' experience of emergency remote teaching
1154 during COVID-19. *Children and Youth Services Review*, 119:105578.
- 1155 [65] Simon, M. A., Tzur, R., Heinz, K., and Kinzel, M. (2004). Explicating a mechanism for
1156 conceptual learning: elaborating the construct of reflective abstraction. *Journal for Research in*
1157 *Mathematics Education*, 35(5):305–329.
- 1158 [66] Snijders, T. A. B. and Bosker, R. (2011). More powerful tests for variance parameters. In
1159 *Multilevel Analysis: An Introduction to Basic and Advanced Multilevel Modeling*, chapter 6, pages
1160 94–108. Sage Publications, 2nd edition.
- 1161 [67] Stepien, K. A. and Baernstein, A. (2006). Education for empathy. *Journal of General Internal*
1162 *Medicine*, 21:524–530.
- 1163 [68] Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B.,
1164 Goyal, N., Hambro, E., Azhar, F., Rodriguz, A., Joulin, A., Grave, E., and Lample, G. (2023).
1165 LLaMA: open and efficient foundation language models. *arXiv*, 2302.13971.
- 1166 [69] Tulchinskii, E., Kuznetsov, K., Kushnareva, L., Cherniavskii, D., Barannikov, S., Pio-
1167 ntkovskaya, I., Nikolenko, S., and Burnaev, E. (2023). Intrinsic dimension estimation for robust
1168 detection of AI-generated texts. *arXiv*, 2306.04723.
- 1169 [70] van Paridon, J., Liu, Q., and Lupyan, G. (2021). How do blind people know that blue is cold?
1170 distributional semantics encode color-adjective associations. *Proceedings of the Annual Meeting of*
1171 *the Cognitive Science Society*, 43(43).

- 1172 [71] Viswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and
1173 Polosukhin, I. (2017). Attention is all you need. In *Advances in Neural Information Processing*
1174 *Systems*.
- 1175 [72] Whalen, J. (2020). Should teachers be trained in emergency remote teaching? Lessons learned
1176 from the COVID-19 pandemic. *Journal of Technology and Teacher Education*, 28(2):189–199.
- 1177 [73] Young, J. R. (2012). Inside the Coursera contract: how an upstart company might profit from
1178 free courses. *The Chronicle of Higher Education*, 19(7):1–4.
- 1179 [74] Ziman, K., Heusser, A. C., Fitzpatrick, P. C., Field, C. E., and Manning, J. R. (2018). Is
1180 automatic speech-to-text transcription ready for use in psychological experiments? *Behavior*
1181 *Research Methods*, 50:2597–2605.