

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 tracking and characterizing the acquisition of conceptual knowledge. Our approach embeds each concept in a high-dimensional representation space, where nearby coordinates reflect similar or related concepts. We test our approach using behavioral data from participants who answered small sets of multiple-choice quiz questions interleaved between watching two course videos from the Khan Academy platform. We apply our framework to the videos' transcripts and the text of the quiz questions to quantify the content of each moment of video and each quiz question. We use these embeddings, along with participants' quiz responses, to track how the learners' knowledge changed after watching each video. Our findings show how a small set of quiz questions may be used to obtain rich and meaningful, high-resolution insights into what 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 imaginary 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 [2, 6, 8, 9, 43]?
46 Or weaving a lecture’s atomic elements (e.g., its component words) into a structured network that
47 describes how those individual elements are related [26]? Conceptual understanding could also
48 involve building a mental model that transcends the meanings of those individual atomic elements
49 by reflecting the deeper meaning underlying the gestalt whole [23, 27, 40].

50 The difference between “understanding” and “memorizing,” as framed by researchers in ed-
51 ucation, cognitive psychology, and cognitive neuroscience (e.g., 14, 16, 19, 27, 40), 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 pas-
55 sage might provide some ability to answer simple questions about it. However, text embedding
56 models (e.g., 3–5, 7, 10, 25, 33) also attempt to capture the deeper meaning *underlying* those atomic
57 elements. These models consider not only the co-occurrences of those elements within and across
58 documents, but also patterns in how those elements appear across different scales (e.g., sentences,
59 paragraphs, chapters, etc.), the temporal and grammatical properties of the elements, and other
60 high-level characteristics of how they are used [28, 29]. According to these models, the deep
61 conceptual meaning of a document may be captured by a feature vector in a high-dimensional
62 representation space, wherein nearby vectors reflect conceptually related documents. A model
63 that succeeds at capturing an analogue of “understanding” is able to assign nearby feature vectors
64 to two conceptually related documents, *even when the specific words contained in those documents have*
65 *very little overlap.*

66 Given these insights, what form might a representation of the sum total of a person’s knowledge
67 take? First, we might require a means of systematically describing or representing the nearly
68 infinite set of possible things a person could know. Second, we might want to account for potential
69 associations between different concepts. For example, the concepts of “fish” and “water” might be
70 associated in the sense that fish live in water. Third, knowledge may have a critical dependency
71 structure, such that knowing about a particular concept might require first knowing about a set of
72 other concepts. For example, understanding the concept of a fish swimming in water first requires

73 understanding what fish and water *are*. Fourth, as we learn, our “current state of knowledge”
74 should change accordingly. Learning new concepts should both update our characterizations of
75 “what is known” and also unlock any now-satisfied dependencies of those newly learned concepts
76 so that they are “tagged” as available for future learning.

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

85 Our primary research goal is to advance our understanding of what it means to acquire deep,
86 real-world conceptual knowledge. Traditional laboratory approaches to studying learning and
87 memory (e.g., list-learning studies) often draw little distinction between memorization and under-
88 standing. Instead, these studies typically focus on whether information is effectively encoded or
89 retrieved, rather than whether the information is *understood*. Approaches to studying conceptual
90 learning, such as category learning experiments, can begin to investigate the distinction between
91 memorization and understanding, often by training participants to distinguish arbitrary or ran-
92 dom features in otherwise meaningless categorized stimuli. However the objective of real-world
93 training, or learning from life experiences more generally, is often to develop new knowledge that
94 may be applied in *useful* ways in the future. In this sense, the gap between modern learning the-
95 ories and modern pedagogical approaches that inform classroom learning strategies is enormous:
96 most of our theories about *how* people learn are inspired by experimental paradigms and models
97 that have only peripheral relevance to the kinds of learning that students and teachers actually
98 seek [16, 27]. To help bridge this gap, our study uses course materials from real online courses to
99 inform, fit, and test models of real-world conceptual learning. We also provide a demonstration of
100 how our models can be used to construct “maps” of what students know, and how their knowl-

101 edge changes with training. In addition to helping to visually capture knowledge (and changes
102 in knowledge), we hope that such maps might lead to real-world tools for improving how we
103 educate. Taken together, our work shows that existing course materials and evaluative tools like
104 short multiple-choice quizzes may be leveraged to gain highly detailed insights into what students
105 know and how they learn.

106 **Results**

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

121 We asked participants in our study to complete brief multiple-choice quizzes before, between,
122 and after watching two lecture videos from the Khan Academy [22] platform (Fig. 1). The first
123 lecture video, entitled *Four Fundamental Forces*, discussed the four fundamental forces in physics:
124 gravity, strong and weak interactions, and electromagnetism. The second, entitled *Birth of Stars*,
125 provided an overview of our current understanding of how stars form. We selected these particular
126 lectures to satisfy three general criteria. First, we wanted both lectures to be accessible to a broad



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.

audience (i.e., with minimal prerequisite knowledge) so as to limit the impact of prior training on our participants' abilities to learn from the lectures. To this end, we selected two introductory videos that were intended to be viewed at the start of students' training in their respective content areas. Second, we wanted both lectures to have some related content, so that we could test our approach's ability to distinguish similar conceptual content. To this end, we chose two videos from the same (per instructor annotations) Khan Academy course domain, "Cosmology and Astronomy." Third, we sought to minimize dependencies and specific overlap between the videos. For example, we did not want participants' abilities to understand one video to (directly) influence their abilities to understand the other. To satisfy this last criterion, we chose videos from two different lecture series (Lectures 1 and 2 were from the "Scale of the Universe" and "Stars, Black Holes, and Galaxies" series, respectively).

We also wrote a set of multiple-choice quiz questions that we hoped would enable us to evaluate participants' knowledge about each individual lecture, along with related knowledge about physics not specifically presented in either video (see Tab. S1 for the full list of questions in our stimulus pool). Participants answered questions randomly drawn from each content area (Lecture 1, Lecture 2, and general physics knowledge) on each of the three quizzes. Quiz 1 was intended to assess participants' "baseline" knowledge before training, Quiz 2 assessed knowledge after watching the *Four Fundamental Forces* video (i.e., Lecture 1), and Quiz 3 assessed knowledge

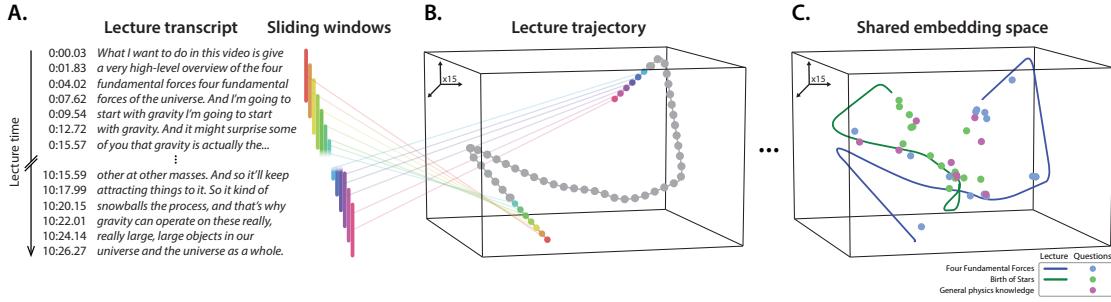


Figure 2: Modeling course content. **A. Building a document pool from sliding windows of text.** We decompose each lecture’s transcript into a series of overlapping sliding windows. The full set of transcript snippets (across all windows) may be treated as a set of “documents” for training a text embedding model. **B. Constructing lecture content trajectories.** After training our 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 (Tab. S1), 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 after watching the *Birth of Stars* video (i.e., Lecture 2).

146 To study in detail how participants’ conceptual knowledge changed over the course of the
 147 experiment, we first sought to model the conceptual content presented to them at each moment
 148 throughout each of the two lectures. We adapted an approach we developed in prior work [17] to
 149 identify the latent themes in the lectures using a topic model [4]. Briefly, topic models take as input
 150 a collection of text documents, and learn a set of “topics” (i.e., latent themes) from their contents.
 151 Once fit, a topic model can be used to transform arbitrary (potentially new) documents into sets
 152 of “topic proportions,” describing the weighted blend of learned topics reflected in their texts. We
 153 parsed automatically generated transcripts of the two lectures into overlapping sliding windows,
 154 where each window contained the text of the lecture transcript from a particular time span. We
 155 treated the set of text snippets (across all of these windows) as documents to fit our model (Fig. 2A;
 156 see *Constructing text embeddings of multiple lectures and questions*). Transforming the text from every
 157 sliding window with our model yielded a number-of-windows by number-of-topics (15) topic-
 158 proportions matrix that described the unique mixture of broad themes from both lectures reflected
 159 in each window’s “topic vector” (i.e., column of the topic-proportions matrix)

160 is analogous to a coordinate in a 15-dimensional space whose axes are topics discovered by the
161 model. Within this space, each lecture’s sequence of topic vectors (i.e., corresponding to its
162 transcript’s overlapping text snippets across sliding windows) forms a *trajectory* that captures how
163 its conceptual content unfolds over time (Fig. 2B). We resampled these trajectories to a resolution
164 of one topic vector for each second of video (i.e., 1 Hz).

165 We hypothesized that a topic model trained on transcripts of the two lectures should also capture
166 the conceptual knowledge probed by each quiz question. If indeed the topic model could capture
167 information about the deeper conceptual content of the lectures (i.e., beyond surface-level details
168 such as particular word choices), then we should be able to recover a correspondence between each
169 lecture and questions *about* each lecture. Importantly, such a correspondence could not solely arise
170 from superficial text matching between lecture transcripts and questions, since the lectures and
171 questions used different words. Simply comparing the average topic weights from each lecture and
172 question set (averaging across time and questions, respectively) reveals a striking correspondence
173 (Fig. S1). Specifically, the average topic weights from Lecture 1 are strongly correlated with the
174 average topic weights from Lecture 1 questions ($r(13) = 0.809, p < 0.001, 95\% \text{ CI} = [0.633, 0.962]$), and the average topic weights from Lecture 2 are strongly correlated with the
175 average topic weights from Lecture 2 questions ($r(13) = 0.728, p = 0.002, 95\% \text{ CI} = [0.456, 0.920]$).
176 At the same time, the average topic weights from the two lectures are *negatively* correlated with
177 their non-matching question sets (Lecture 1 video vs. Lecture 2 questions: $r(13) = -0.547, p = 0.035,$
178 $95\% \text{ CI} = [-0.812, -0.231]$; Lecture 2 video vs. Lecture 1 questions: $r(13) = -0.612, p = 0.015, 95\%$
179 $\text{CI} = [-0.874, -0.281]$), indicating that the topic model also exhibits some degree of specificity. The
180 full set of pairwise comparisons between average topic weights for the lectures and question sets
181 is reported in Figure S1.

183 Another, more sensitive, way of summarizing the conceptual content of the lectures and ques-
184 tions is to look at *variability* in how topics are weighted over time and across different questions
185 (Fig. 3). Intuitively, the variability in the expression of a given topic relates to how much “infor-
186 mation” [13] the lecture (or question set) reflects about that topic. For example, suppose a given
187 topic is weighted on heavily throughout a lecture. That topic might be characteristic of some



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 Table S2. **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 aspect or property of the lecture *overall* (conceptual or otherwise), but unless the topic’s weights
 189 changed in meaningful ways over time, the topic would be a poor indicator of any *specific* concep-
 190 tual content in the lecture. We therefore also compared the variances in topic weights (across time
 191 or questions) between the lectures and questions. The variability in topic expression (over time
 192 and across questions) was similar for the Lecture 1 video and questions ($r(13) = 0.824, p < 0.001,$
 193 $95\% \text{ CI} = [0.696, 0.973]$) and the Lecture 2 video and questions ($r(13) = 0.801, p < 0.001, 95\%$
 194 $\text{CI} = [0.539, 0.958]$). However, as reported in Figure 3B, the variability in topic expressions across
 195 *different* videos and lecture-specific questions (i.e., Lecture 1 video vs. Lecture 2 questions; Lecture 2
 196 video vs. Lecture 1 questions) were negatively correlated, and neither video’s topic variability was
 197 reliably correlated with the topic variability across general physics knowledge questions. Taken
 198 together, the analyses reported in Figures 3 and S1 indicate that a topic model fit to the videos’
 199 transcripts can also reveal correspondences (at a coarse scale) between the lectures and questions.
 200 Although a single lecture may be organized around a single broad theme at a coarse scale, at a
 201 finer scale each moment of a lecture typically covers a narrower range of content. We wondered

202 whether a text embedding model trained on the lectures' transcripts might capture some of this
203 finer scale content. For example, if a particular question asks about the content from one small part
204 of a lecture, we wondered whether the text embeddings could be used to automatically identify
205 the "matching" moment(s) in the lecture. When we correlated each question's topic vector with
206 the topic vectors from each second of the lectures, we found some evidence that each question is
207 temporally specific (Fig. 4). In particular, most questions' topic vectors were maximally correlated
208 with a well-defined (and relatively narrow) range of timepoints from their corresponding lectures,
209 and the correlations fell off sharply outside of that range. We also qualitatively examined the best-
210 matching intervals for each question by comparing the text of the question to the text of the most-
211 correlated parts of the lectures. Despite that the questions were excluded from the text embedding
212 model's training set, in general we found (through manual inspection) a close correspondence
213 between the conceptual content that each question probed and the content covered by the best-
214 matching moments of the lectures. Two representative examples are shown at the bottom of
215 Figure 4.

216 The ability to quantify how much each question is "asking about" the content from each moment
217 of the lectures could enable high-resolution insights into participants' knowledge. Traditional
218 approaches to estimating how much a student "knows" about the content of a given lecture entail
219 computing the proportion of correctly answered questions. But if two students receive identical
220 scores on an exam, might our modeling framework help us to gain more nuanced insights into
221 the *specific* content that each student has mastered (or failed to master)? For example, a student
222 who misses three questions that were all about the same concept (e.g., concept *A*) will have gotten
223 the same *proportion* of questions correct as another student who missed three questions about
224 three *different* concepts (e.g., *A*, *B*, and *C*). But if we wanted to fill in the "gaps" in the two
225 students' understandings, we might do well to focus on concept *A* for the first student, but to
226 also add in materials pertaining to concepts *B* and *C* for the second student. In other words, raw
227 "proportion-correct" measures may capture *how much* a student knows, but not *what* they know.
228 We wondered whether our modeling framework might enable us to (formally and automatically)
229 infer participants' knowledge at the scale of individual concepts (e.g., as captured by a single

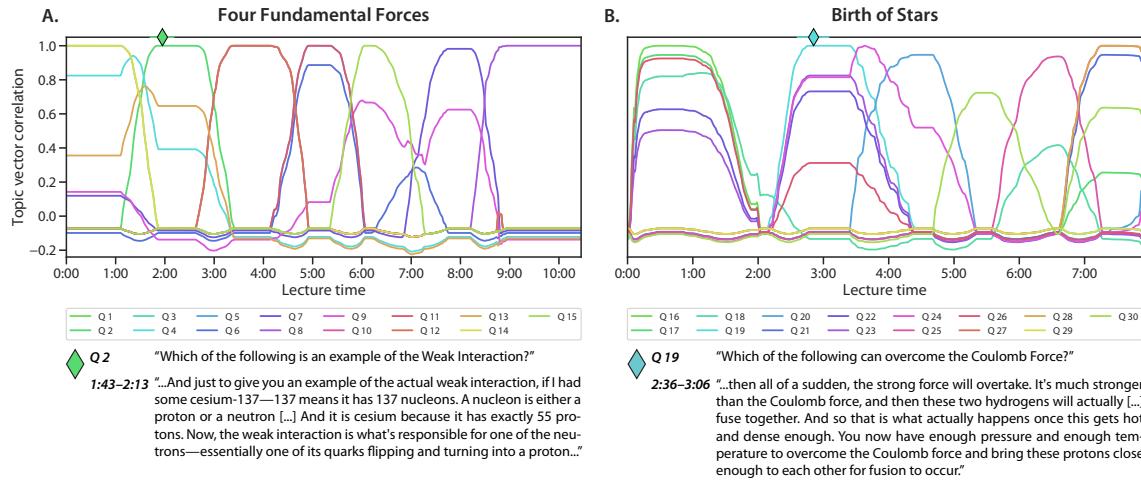


Figure 4: Which parts of each lecture are captured by each question? Each panel displays timeseries plots showing how each question’s topic vector correlates with each video timepoint’s topic vector (Panel A.: correlations for the *Four Fundamental Forces* lecture and associated questions; Panel B.: 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.

230 moment of a lecture).

231 We developed a simple formula (Eqn. 1) for using a participant’s responses to a small set
232 of multiple-choice questions to estimate how much the participant “knows” about the concept
233 reflected by any arbitrary coordinate, x , in text embedding space (e.g., the content reflected by
234 any moment in a lecture they had watched; see *Estimating dynamic knowledge traces*). Essentially,
235 the estimated knowledge at the coordinate is given by the weighted average proportion of quiz
236 questions the participant answered correctly, where the weights reflect how much each question
237 is “about” the content at x . When we apply this approach to estimate the participant’s knowledge
238 about the content presented in each moment of each lecture, we can obtain a detailed timecourse
239 describing how much “knowledge” the participant has about any part of the lecture. As shown
240 in Figure 5, we can also apply this approach separately for the questions from each quiz the
241 participants took throughout the experiment. From just a few questions per quiz (see Sec.), we
242 obtain a high-resolution snapshot (at the time each quiz was taken) of what the participants knew
243 about any moment’s content, from either of the two lectures they watched (comprising a total of
244 1,100 samples across the two lectures).

245 Of course, even though the timecourses in Figure 5A and C provide detailed *estimates* about
246 participants’ knowlege, those estimates are only *useful* to the extent that they accurately reflect what
247 participants actually know. As one sanity check, we anticipated that the knowledge estimates
248 should show a content-specific “boost” in participants’ knowledge after watching each lecture.
249 In other words, if participants learn about each lecture’s content when they watch each lecture,
250 the knowledge estimates should reflect that. After watching the *Four Fundamental Forces* lecture,
251 participants should show more knowledge for the content of that lecture than they had before,
252 and that knowledge should persist for the remainder of the experiment. Specifically, knowledge
253 about that lecture’s content should be relatively low when estimated using Quiz 1 responses,
254 but should increase when estimated using Quiz 2 or 3 responses (Fig. 5B). Indeed, we found
255 that participants’ estimated knowledge about the content of the *Four Fundamental Forces* was
256 substantially higher on Quiz 2 versus Quiz 1 ($t(49) = 8.764$, $p < 0.001$) and on Quiz 3 versus
257 Quiz 1 ($t(49) = 10.519$, $p < 0.001$). We found no reliable differences in estimated knowledge about



Figure 5: Estimating moment-by-moment knowledge acquisition. **A. Moment-by-moment knowledge about the *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 one quiz (color). The traces are averaged across participants. **B. Average estimated knowledge about the *Four Fundamental Forces*.** Each bar displays the across-timepoint average knowledge, estimated using the responses to one quiz's questions. **C. Moment-by-moment knowledge about the *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 the *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.

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

If we are able to accurately estimate a participant's knowledge about the content tested by a given question, our estimates of their knowledge should carry some predictive information about whether the participant is likely to answer the question correctly or incorrectly. We developed a statistical approach to test this claim. For each question in turn, for each participant, we used Equation 1 to estimate (using all *other* questions from the same quiz, from the same participant) the participant's knowledge at the held-out question's embedding coordinate. For each quiz, we grouped these estimates into two distributions: one for the estimated knowledge at the coordinates of each *correctly* answered question, and another for the estimated knowledge at the coordinates of each *incorrectly* answered question (Fig. 6). We then used independent samples *t*-tests to compare the means of these distributions of estimated knowledge.

For the initial quizzes participants took (prior to watching either lecture), participants' estimated knowledge tended to be low overall, and relatively unstructured (Fig. 6, left panel). When we held out individual questions and estimated their knowledge at the held-out questions' embedding coordinates, we found no reliable differences in the estimates when the held-out question had been correctly versus incorrectly answered ($t(633) = 0.577, p = 0.564$). After watching the first video, estimated knowledge for held-out correctly answered questions (from the second quiz; Fig. 6, middle panel) exhibited a positive shift relative to held-out incorrectly answered questions ($t(633) = 3.961, p < 0.001$). After watching the second video, estimated knowledge (from the third quiz; Fig. 6, right panel) for *all* questions exhibited a positive shift. However, the increase

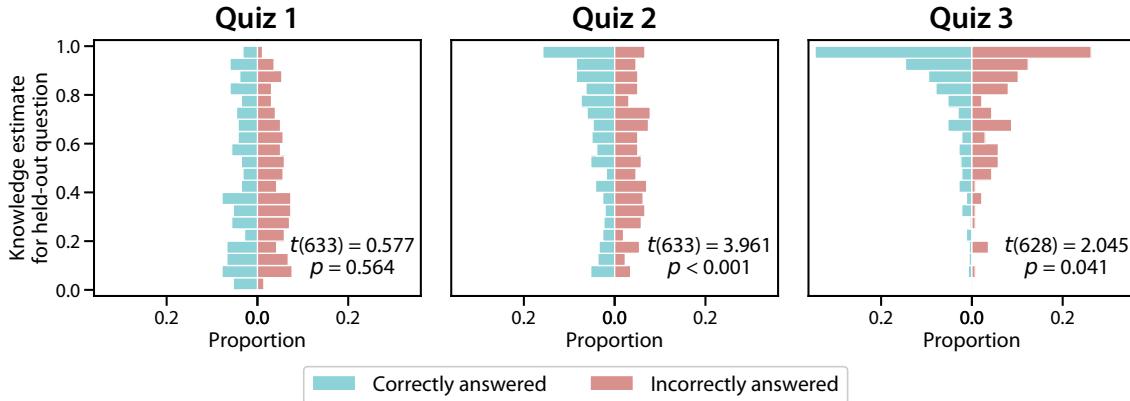


Figure 6: Estimating knowledge at the embedding coordinates of held-out questions. Separately for each quiz (panel), we plot the distributions of predicted knowledge at the embedding coordinates of each held-out correctly (blue) or incorrectly (red) answered question. The t -tests reported in each panel are between the distributions of estimated knowledge at the coordinates of correctly versus incorrectly answered held-out questions.

in estimated knowledge for held-out correctly answered questions was larger than for held-out incorrectly answered questions (estimated knowledge for correctly versus incorrectly answered Quiz 3 questions: $t(628) = 2.045$, $p = 0.041$).

Knowledge estimates need not be limited to the content of the lectures. As illustrated in Figure 7, our general approach to estimating knowledge from a small number of quiz questions may be applied to *any* content, given its text embedding coordinate. To visualize how knowledge “spreads” through text embedding space to content beyond the lectures participants watched, we first fit a new topic model to the lectures’ sliding windows with $k = 100$ topics. We hoped that increasing the number of topics from 15 to 100 might help us to generalize the knowledge predictions. (Aside from increasing the number of topics from 15 to 100, all other procedures and model parameters were carried over from the preceding analyses.) As in our other analyses, we resampled each lecture’s topic trajectory to 1 Hz and also projected each question into a shared text embedding space.

We projected the resulting 100-dimensional topic vectors (for each second of video and for each question) onto a shared 2-dimensional plane (see *Creating knowledge and learning map visualizations*). Next, we sampled points from a 100×100 grid of coordinates that evenly tiled a rectangle enclosing



Figure 7: 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 each lecture are indicated by dotted lines, and the coordinates of each question are indicated by dots. Each map reflects an average across all participants. For individual participants’ maps, see Figures S2, S3, and S4.

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 Figures S5 and S6.

C. Word clouds for sampled points in topic space. Each word cloud displays the relative weights of each word (via their relative sizes) reflected by the blend of topics represented at the locations of the stars on the maps. The words’ colors indicate how much each word is weighted, on average, across all timepoints’ topic vectors in the *Four Fundamental Forces* (blue) and *Birth of Stars* (green) videos, respectively.

302 the 2D projections of the videos and questions. We used Equation 4 to estimate participants' knowl-
303 edge at each of these 10,000 sampled locations, and averaged these estimates across participants to
304 obtain an estimated average *knowledge map* (Fig. 7A). Intuitively, the knowledge map constructed
305 from a given quiz's responses provides a visualization of how "much" participants know about
306 any content expressible by the fitted text embedding model.

307 Several features of the resulting knowledge maps are worth noting. The average knowledge
308 map estimated from Quiz 1 responses (Fig. 7A, leftmost map) shows that participants tended to
309 have relatively little knowledge about any parts of the text embedding space (i.e., the shading is
310 relatively dark everywhere). The knowledge map estimated from Quiz 2 responses shows a marked
311 increase in knowledge on the left side of the map (around roughly the same range of coordinates
312 traversed by the *Four Fundamental Forces* lecture, indicated by the dotted blue line). In other words,
313 participants' estimated increase in knowledge is localized to conceptual content that is nearby (i.e.,
314 related to) the content from the lecture they watched prior to taking Quiz 2. This localization
315 is non-trivial: the knowledge estimates are informed only by the embedded coordinates of the
316 *quiz questions*, not by the embeddings of either lecture (see Eqn. 4). Finally, the knowledge map
317 estimated from Quiz 3 responses shows a second increase in knowledge, localized to the region
318 surrounding the embedding of the *Birth of Stars* lecture participants watched immediately prior to
319 taking Quiz 3.

320 Another way of visualizing these content-specific increases in knowledge after participants
321 viewed each lecture is displayed in Figure 7B. Taking the point-by-point difference between the
322 knowledge maps estimated from responses to a successive pair of quizzes yields a *learning map*
323 that describes the *change* in knowledge estimates from one quiz to the next. These learning maps
324 highlight that the estimated knowledge increases we observed across maps were specific to the
325 regions around the embeddings of each lecture in turn.

326 Because the 2D projection we used to construct the knowledge and learning maps is invertible,
327 we may gain additional insights into these maps' meaning by reconstructing the original high-
328 dimensional topic vector for any location on the map we are interested in. For example, this could
329 serve as a useful tool for an instructor looking to better understand which content areas a student

330 (or a group of students) knows well (or poorly). As a demonstration, we show the top-weighted
331 words from the blends of topics reconstructed from three example locations on the maps (Fig. 7C):
332 one point near the *Four Fundamental Forces* embedding (yellow); a second point near the *Birth of*
333 *Stars* embedding (orange), and a third point between the two lectures' embeddings (pink). As
334 shown in the word clouds in the Panel, the top-weighted words at the example coordinate near
335 the *Four Fundamental Forces* embedding also tended to be weighted heavily by the topics expressed
336 in that lecture. Similarly, the top-weighted words at the example coordinate near the *Birth of Stars*
337 embedding tended to be weighted most heavily by the topics expressed in *that* lecture. And the
338 top-weighted words at the example coordinate between the two lectures' embeddings show a
339 roughly even mix of words most strongly associated with each lecture.

340 Discussion

341 We developed a computational framework that uses short multiple-choice quizzes to gain nuanced
342 insights into what learners know and how their knowledge changes with training. First, we show
343 that our approach can automatically match the conceptual knowledge probed by individual quiz
344 questions to the corresponding moments in lecture videos when those concepts were presented
345 (Fig. 4). Next, we demonstrate how we can estimate moment-by-moment “knowledge traces”
346 that reflect the degree of knowledge participants have about each video’s time-varying content,
347 and capture temporally specific increases in knowledge after viewing each lecture (Fig. 5). We
348 also show that these knowledge estimates can generalize to held-out questions (Fig. 6). Finally,
349 we use our framework to construct visual maps that provide snapshot estimates of how much
350 participants know about any concept within the scope of our text embedding model, and how
351 much their knowledge changes with training (Fig. 7).

352 Over the past several years, the global pandemic has forced many educators to teach re-
353 motely [21, 34, 42, 45]. This change in world circumstances is happening alongside (and perhaps
354 accelerating) geometric growth in the availability of high quality online courses on platforms such
355 as Khan Academy [22], Coursera [46], EdX [24], and others [39]. Continued expansion of the global

356 internet backbone and improvements in computing hardware have also facilitated improvements
357 in video streaming, enabling videos to be easily shared and viewed by large segments of the
358 world's population. This exciting time for online course instruction provides an opportunity to
359 re-evaluate how we, as a global community, educate ourselves and each other. For example, we
360 can ask: what makes an effective course or training program? Which aspects of teaching might
361 be optimized and/or augmented by automated tools? How and why do learning needs and goals
362 vary across people? How might we lower barriers to achieving a high-quality education?

363 Alongside these questions, there is a growing desire to extend existing theories beyond the
364 domain of lab testing rooms and into real classrooms [20]. In part, this has led to a recent
365 resurgence of "naturalistic" or "observational" experimental paradigms that attempt to better
366 reflect more ethologically valid phenomena that are more directly relevant to real-world situations
367 and behaviors [35]. In turn, this has brought new challenges in data analysis and interpretation. A
368 key step towards solving these challenges will be to build explicit models of real-world scenarios
369 and how people behave in them (e.g., models of how people learn conceptual content from real-
370 world courses, as in our current study). A second key step will be to understand which sorts of
371 signals derived from behaviors and/or other measurements (e.g., neurophysiological data; 1, 12, 32,
372 36, 37) might help to inform these models. A third major step will be to develop and employ reliable
373 ways of evaluating the complex models and data that are a hallmark of naturalistic paradigms.

374 Beyond specifically predicting what people *know*, the fundamental ideas we develop here also
375 relate to the notion of "theory of mind" of other individuals [15, 18, 31]. Considering others' unique
376 perspectives, prior experiences, knowledge, goals, etc., can help us to more effectively interact and
377 communicate [38, 41, 44]. One could imagine future extensions of our work (e.g., analogous to
378 the knowledge and learning maps shown in Fig. 7), that attempt to characterize how well-aligned
379 different people's knowledge bases or backgrounds are. In turn, this might be used to model how
380 knowledge (or other forms of communicable information) flows not just between teachers and
381 students, but between friends having a conversation, individuals on a first date, participants at
382 a business meeting, doctors and patients, experts and non-experts, political allies or adversaries,
383 and more. For example, the extent to which two people's knowledge maps "match" or "align" in

384 a given region of text embedding space might serve as a predictor of how effectively they will be
385 able to communicate about the corresponding conceptual content.

386 Ultimately, our work suggests a rich new line of questions about the geometric “form” of
387 knowledge, how knowledge changes over time, and how we might map out the full space of
388 what an individual knows. Our finding that detailed estimates about knowledge may be obtained
389 from short quizzes shows one way that traditional approaches to evaluation in education may be
390 extended. We hope that these advances might help pave the way for new approaches to teaching
391 or delivering educational content that are tailored to individual students’ learning needs and goals.

392 Materials and methods

393 Participants

394 We enrolled a total of 50 Dartmouth undergraduate students in our study. Participants received
395 course credit for enrolling. We asked each participant to complete a demographic survey that
396 included questions about their age, gender, native spoken language, ethnicity, race, hearing, color
397 vision, sleep, coffee consumption, level of alertness, and several aspects of their educational back-
398 ground and prior coursework.

399 Participants’ ages ranged from 18 to 22 years (mean: 19.52 years; standard deviation: 1.09
400 years). A total of 15 participants reported their gender as male and 35 participants reported their
401 gender as female. A total of 49 participants reported their native language as “English” and 1
402 reported having another native language. A total of 47 participants reported their ethnicity as
403 “Not Hispanic or Latino” and three reported their ethnicity as “Hispanic or Latino.” Participants
404 reported their races as White (32 participants), Asian (14 participants), Black or African American
405 (5 participants), American Indian or Alaska Native (1 participant), and Native Hawaiian or Other
406 Pacific Islander (1 participant). (Note that some participants selected multiple racial categories.)

407 A total of 49 participants reporting having normal hearing and 1 participant reported having
408 some hearing impairment. A total of 49 participants reported having normal color vision and 1

409 participant reported being color blind. Participants reported having had, on the night prior to
410 testing, 2–4 hours of sleep (1 participant), 4–6 hours of sleep (9 participants), 6–8 hours of sleep (35
411 participants), or 8+ hours of sleep (5 participants). They reported having consumed, on the same
412 day and leading up to their testing session, 0 cups of coffee (38 participants), 1 cup of coffee (10
413 participants), 3 cups of coffee (1 participant), or 4+ cups of coffee (1 participant).

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

419 Participants reported their undergraduate major(s) as “social sciences” (28 participants), “nat-
420 ural sciences” (16 participants), “professional” (e.g., pre-med or pre-law; 8 participants), “mathe-
421 matics and engineering” (7 participants), “humanities” (4 participants), or “undecided” (3 partici-
422 pants). Note that some participants selected multiple categories for their undergraduate major(s).
423 We also asked participants about the courses they had taken. In total, 45 participants reported hav-
424 ing taken at least one Khan Academy course in the past, and 5 reported not having taken any Khan
425 Academy courses. Of those who reported having watched at least one Khan Academy course,
426 7 participants reported having watched 1–2 courses, 11 reported having watched 3–5 courses, 8
427 reported having watched 5–10 courses, and 19 reported having watched 10 or more courses. We
428 also asked participants about the specific courses they had watched, categorized under different
429 subject areas. In the “Mathematics” area, participants reported having watched videos on AP
430 Calculus AB (21 participants), Precalculus (17 participants), Algebra 2 (14 participants), AP Cal-
431 culus BC (12 participants), Trigonometry (11 participants), Algebra 1 (10 participants), Geometry
432 (8 participants), Pre-algebra (7 participants), Multivariable Calculus (5 participants), Differential
433 Equations (5 participants), Statistics and Probability (4 participants), AP Statistics (2 participants),
434 Linear Algebra (2 participants), Early Math (1 participant), Arithmetic (1 participant), and other
435 videos not listed in our survey (5 participants). In the “Science and engineering” area, participants
436 reported having watched videos on Chemistry, AP Chemistry, or Organic Chemistry (21 partic-

ipants); Physics, AP Physics I, or AP Physics II (18 participants); Biology, AP Biology; or High school Biology (15 participants); Health and Medicine (1 participant); or other videos not listed in our survey (5 participants). We also asked participants whether they had specifically seen the videos used in our experiment. Of the 45 participants who reported having taken at least one Khan Academy course in the past, 44 participants reported that they had not watched the *Four Fundamental Forces* video, and 1 participant reported that they were not sure whether they had watched it. All participants reported that they had not watched the *Birth of Stars* video. When we asked participants about non-Khan Academy online courses, they reported having watched or taken courses on Mathematics (15 participants), Science and engineering (11 participants), Test preparation (9 participants), Economics and finance (3 participants), Arts and humanities (2 participants), Computing (2 participants), and other categories not listed in our survey (17 participants). Finally, we asked participants about in-person courses they had taken in different subject areas. They reported taking courses in Mathematics (38 participants), Science and engineering (37 participants), Arts and humanities (34 participants), Test preparation (27 participants), Economics and finance (26 participants), Computing (14 participants), College and careers (7 participants), or other courses not listed in our survey (6 participants).

Experiment

We hand-selected two course videos from the Khan Academy platform: *Four Fundamental Forces* (an introduction to gravity, electromagnetism, the weak nuclear force, and the strong nuclear force; duration: 10 minutes and 29 seconds) and *Birth of Stars* (an introduction to how stars are formed; duration: 7 minutes and 57 seconds). We then hand-created 39 multiple-choice questions: 15 about the conceptual content of *Four Fundamental Forces* (i.e., Lecture 1), 15 about the conceptual content of *Birth of Stars* (i.e., Lecture 2), and 9 questions that tested for general conceptual knowledge about basic physics (covering material that was not presented in either video). The full set of questions and answer choices may be found in Table S1.

Over the course of the experiment, participants completed three 13-question multiple-choice quizzes: the first before viewing Lecture 1, the second between Lectures 1 and 2, and the third

464 after viewing Lecture 2 (Fig. 1). The questions appearing on each quiz, for each participant, were
465 randomly chosen from the full set of 39, with the constraints that (a) each quiz contain 5 questions
466 about Lecture 1, 5 questions about Lecture 2, and 3 questions about general physics knowledge,
467 and (b) each question appear exactly once for each participant. The orders of questions on each
468 quiz, and the orders of answer options for each question, were also randomized. Our experimental
469 protocol was approved by the Committee for the Protection of Human Subjects at Dartmouth
470 College. We used the experiment to develop and test our computational framework for estimating
471 knowledge and learning.

472 **Analysis**

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

474 We adapted an approach we developed in prior work [17] to embed each moment of the two
475 lectures and each question in our pool in a common representational space. Briefly, our approach
476 uses a topic model (Latent Dirichlet Allocation; 4), trained on a set of documents, to discover a set
477 of k “topics” or “themes.” Formally, each topic is defined as a set of weights over each word in
478 the model’s vocabulary (i.e., the union of all unique words, across all documents, excluding “stop
479 words.”). Conceptually, each topic is intended to give larger weights to words that are semantically
480 related or tend to co-occur in the same documents. After fitting a topic model, each document
481 in the training set, or any *new* document that contains at least some of the words in the model’s
482 vocabulary, may be represented as a k -dimensional vector describing how much the document
483 (most probably) reflects each topic. (Unless, otherwise noted, we used $k = 15$ topics.)

484 As illustrated in Figure 2A, we start by building up a corpus of documents using overlapping
485 sliding windows that span each video’s transcript. Khan Academy provides professionally created,
486 manual transcriptions of all videos for closed captioning. However, such transcripts would not
487 be readily available in all contexts to which our framework could potentially be applied. Khan
488 Academy videos are hosted on the YouTube platform, which additionally provides automated
489 captions. We opted to use these automated transcripts (which, in prior work, we have found are

490 of sufficiently near-human quality yield reliable data in behavioral studies; 47) when developing
491 our framework in order to make it more directly extensible and adaptable by others in the future.

492 We fetched these automated transcripts using the `youtube-transcript-api` Python pack-
493 age [11]. The transcripts consisted of one timestamped line of text for every few seconds (mean:
494 2.34 s; standard deviation: 0.83 s) of spoken content in the video (i.e., corresponding to each indi-
495 vidual caption that would appear on-screen if viewing the lecture via YouTube, and when those
496 lines would appear). We defined a sliding window length of (up to) $w = 30$ transcript lines, and
497 assigned each window a timestamp corresponding to the midpoint between its first and last lines'
498 timestamps. These sliding windows ramped up and down in length at the very beginning and
499 end of the transcript, respectively. In other words, the first sliding window covered only the first
500 line from the transcript; the second sliding window covered the first two lines; and so on. This
501 insured that each line of the transcript appeared in the same number (w) of sliding windows. After
502 performing various standard text preprocessing (e.g., normalizing case, lemmatizing, removing
503 punctuation and stop-words), we treated the text from each sliding window as a single “docu-
504 ment,” and combined these documents across the two videos’ windows to create a single training
505 corpus for the topic model. The top words from each of the 15 discovered topics may be found in
506 Table S2.

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

514 We transformed each sliding window’s text into a topic vector, and then used linear interpo-
515 lation (independently for each topic dimension) to resample the resulting timeseries to one vector
516 per second. We also used the fitted model to obtain topic vectors for each question in our pool
517 (Tab. S1). Taken together, we obtained a *trajectory* for each video, describing its path through topic

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

522 **Estimating dynamic knowledge traces**

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

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

525 where

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

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

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

540 **Creating knowledge and learning map visualizations**

541 An important feature of our approach is that, given a trained text embedding model and partic-
542 ipants’ quiz performance on each question, we can estimate their knowledge about *any* content
543 expressible by the embedding model—not solely the content explicitly probed by the quiz ques-
544 tions or even appearing in the lectures. To visualize these estimates (Figs. 7, S2, S3, S4, S5, and S6),
545 we used Uniform Manifold Approximation and Projection (UMAP; 30) to construct a 2D projection
546 of the text embedding space. Sampling the original 100-dimensional space at high resolution to
547 obtain an adequate set of topic vectors spanning the embedding space would be computationally
548 intractable. However, sampling a 2D grid is trivial.

549 At a high level, the UMAP algorithm obtains low-dimensional embeddings by minimizing
550 the cross-entropy between the pairwise (clustered) distances between the observations in their
551 original (e.g., 100-dimensional) space and the pairwise (clustered) distances in the low-dimensional
552 embedding space (in our approach, the embedding space is 2D). In our implementation, pairwise
553 distances in the original high-dimensional space were defined as 1 minus the correlation between
554 the pair of coordinates, and pairwise distances in the low-dimensional embedding space were
555 defined as the Euclidean distance between the pair of coordinates.

556 In our application, all of the coordinates we embedded were topic vectors, whose elements
557 are always non-negative. Although UMAP is an invertible transformation at the embedding
558 locations of the original data, other locations in the embedding space will not necessarily follow
559 the same implicit “rules” as the original high-dimensional data. For example, inverting an arbitrary
560 coordinate in the embedding space might result in negative-valued vectors, which are incompatible
561 with the topic modeling framework. To protect against this issue, we log-transformed the topic
562 vectors prior to embedding them in the 2D space. When we inverted the embedded vectors (e.g.,
563 to estimate topic vectors or word clouds, as in Fig. 7C), we passed the inverted (log-transformed)
564 values through the exponential function to obtain a vector of non-negative values.

565 After embedding both lectures’ topic trajectories and the topic vectors of every question, we
566 defined a rectangle enclosing the 2D projections of the lectures’ and quizzes’ embeddings. We then

567 sampled points from a regular 100×100 grid of coordinates that evenly tiled this enclosing rectangle.
568 We sought to estimate participants' knowledge (and learning, i.e., changes in knowledge) at each
569 of the resulting 10,000 coordinates.

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

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

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

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

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

581 Author contributions

582 Conceptualization: PCF, ACH, and JRM. Methodology: PCF, ACH, and JRM. Software: PCF.
583 Validation: PCF. Formal analysis: PCF. Resources: PCF, ACH, and JRM. Data curation: PCF.
584 Writing (original draft): JRM. Writing (review and editing): PCF, ACH, and JRM. Visualization:
585 PCF and JRM. Supervision: JRM. Project administration: PCF. Funding acquisition: JRM.

586 **Data and code availability**

587 All of the data analyzed in this manuscript, along with all of the code for running our experiment
588 and carrying out the analyses may be found at [https://github.com/ContextLab/efficient-learning-](https://github.com/ContextLab/efficient-learning-khan)
589 [khan](#).

590 **Acknowledgements**

591 We acknowledge useful discussions, assistance in setting up an earlier (unpublished) version of
592 this study, and assistance with data collection efforts from Will Baxley, Max Bluestone, Daniel
593 Carstensen, Kunal Jha, Caroline Lee, Lucy Owen, Xinming Xu, and Kirsten Ziman. Our work
594 was supported in part by NSF CAREER Award Number 2145172 to JRM. The content is solely the
595 responsibility of the authors and does not necessarily represent the official views of our supporting
596 organizations. The funders had no role in study design, data collection and analysis, decision to
597 publish, or preparation of the manuscript.

598 **References**

- 599 [1] Bevilacque, D., Davidesco, I., Wan, L., and Chaloner, K. (2019). Brain-to-brain synchrony and
600 learning outcomes vary by student-teacher dynamics: evidence from a real-world classroom
601 electroencephalography study. *Journal of Cognitive Neuroscience*, 31(3):401–411.
- 602 [2] Blaye, A., Bernard-Peyron, V., Paour, J.-L., and Bonthoux, F. (2006). Category flexibility in chil-
603 dren: distinguishing response flexibility from conceptual flexibility; the protracted development
604 of taxonomic representations. *European Journal of Developmental Psychology*, 3(2):163–188.
- 605 [3] Blei, D. M. and Lafferty, J. D. (2006). Dynamic topic models. In *Proceedings of the International*
606 *Conference on Machine Learning*, pages 113–120, New York, NY. Association for Computing
607 Machinery.

- 608 [4] Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine*
609 *Learning Research*, 3:993–1022.
- 610 [5] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A.,
611 Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child,
612 R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M.,
613 Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei,
614 D. (2020). Language models are few-shot learners. *arXiv*, 2005.14165.
- 615 [6] Caramazza, A. and Mahon, B. Z. (2003). The organization of conceptual knowledge: the
616 evidence from category-specific semantic deficits. *Trends in Cognitive Sciences*, 7(8):354–361.
- 617 [7] Cer, D., Yang, Y., Kong, S. Y., Hua, N., Limtiaco, N., John, R. S., Constant, N., Guajardo-
618 Cespedes, M., Yuan, S., Tar, C., Sung, Y.-H., Strope, B., and Kurzweil, R. (2018). Universal
619 sentence encoder. *arXiv*, 1803.11175.
- 620 [8] Constantinescu, A. O., O'Reilly, J. X., and Behrens, T. E. J. (2016). Organizing conceptual
621 knowledge in humans with a gridlike code. *Science*, 352(6292):1464–1468.
- 622 [9] Deacon, D., Grose-Fifer, J., Yang, C. M., Stanick, V., Hewitt, S., and Dynowska, A. (2004).
623 Evidence for a new conceptualization of semantic representation in the left and right cerebral
624 hemispheres. *Cortex*, 40(3):467–478.
- 625 [10] Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., and Harshman, R. (1990).
626 Indexing by latent semantic analysis. *Journal of the American Society for Information Science*,
627 41(6):391–407.
- 628 [11] Depoix, J. (2019). YouTube transcript/subtitle API. <https://github.com/jdepoix/youtube-transcript-api>.
- 629
- 630 [12] Dikker, S., Wan, L., Davidesco, I., Kaggen, L., Oostrik, M., McClintock, J., Rowland, J.,
631 Michalareas, G., van Bavel, J. J., Ding, M., and Poeppel, D. (2017). Brain-to-brain synchrony
632 tracks real-world dynamic group interactions in the classroom. *Current Biology*, 27(9):1375–1380.

- 633 [13] Fisher, R. A. (1922). On the mathematical foundations of theoretical statistics. *Philosophical*
634 *Transactions of the Royal Society A*, 222(602):309–368.
- 635 [14] Gallagher, J. J. (2000). Teaching for understanding and application of science knowledge.
636 *School Science and Mathematics*, 100(6):310–318.
- 637 [15] Goldstein, T. R. and Winner, E. (2012). Enhancing empathy and theory of mind. *Journal of*
638 *Cognition and Development*, 13(1):19–37.
- 639 [16] Hall, R. and Greeno, J. (2008). *21st century education: A reference handbook*, chapter Conceptual
640 learning, pages 212–221. Sage Publications.
- 641 [17] Heusser, A. C., Fitzpatrick, P. C., and Manning, J. R. (2021). Geometric models reveal be-
642 havioral and neural signatures of transforming naturalistic experiences into episodic memories.
643 *Nature Human Behavior*, 5:905–919.
- 644 [18] Kanske, P., Böckler, A., and Singer, T. (2015). Models, mechanisms and moderators dissociating
645 empathy and theory of mind. In *Social Behavior from Rodents to Humans*, pages 193–206. Springer.
- 646 [19] Katona, G. (1940). *Organizing and memorizing: studies in the psychology of learning and teaching*.
647 Columbia University Press.
- 648 [20] Kaufman, D. M. (2003). Applying educational theory in practice. *British Medical Journal*,
649 326(7382):213–216.
- 650 [21] Kawasaki, H., Yamasaki, S., Masuoka, Y., Iwasa, M., Fukita, S., and Matsuyama, R. (2021).
651 Remote teaching due to COVID-19: an exploration of its effectiveness and issues. *International*
652 *Journal of Environmental Research and Public Health*, 18(5):2672.
- 653 [22] Khan, S. (2004). *The Khan Academy*. Salman Khan.
- 654 [23] Kintsch (1970). *Learning, memory, and conceptual processes*. Wiley.
- 655 [24] Kolowich, S. (2013). How EdX plans to earn, and share, revenue from its free online courses.
656 *The Chronicle of Higher Education*, 21:1–5.

- 657 [25] Landauer, T. K. and Dumais, S. T. (1997). A solution to Plato's problem: the latent semantic
658 analysis theory of acquisition, induction, and representation of knowledge. *Psychological Review*,
659 104:211–240.
- 660 [26] Lee, H. and Chen, J. (2022). Predicting memory from the network structure of naturalistic
661 events. *Nature Communications*, 13(4235):doi.org/10.1038/s41467-022-31965-2.
- 662 [27] Maclellan, E. (2005). Conceptual learning: the priority for higher education. *British Journal of
663 Educational Studies*, 53(2):129–147.
- 664 [28] Manning, J. R. (2020). Context reinstatement. In Kahana, M. J. and Wagner, A. D., editors,
665 *Handbook of Human Memory*. Oxford University Press.
- 666 [29] Manning, J. R. (2021). Episodic memory: mental time travel or a quantum “memory wave”
667 function? *Psychological Review*, 128(4):711–725.
- 668 [30] McInnes, L., Healy, J., and Melville, J. (2018). UMAP: Uniform manifold approximation and
669 projection for dimension reduction. *arXiv*, 1802(03426).
- 670 [31] Meltzoff, A. N. (2011). Social cognition and the origins of imitation, empathy, and theory of
671 mind. In *The Wiley-Blackwell handbook of childhood cognitive development*. Wiley-Blackwell.
- 672 [32] Meshulam, M., Hasenfratz, L., Hillman, H., Liu, Y. F., Nguyen, M., Norman, K. A., and Hasson,
673 U. (2020). Neural alignment predicts learning outcomes in students taking an introduction to
674 computer science course. *Nature Communications*, 12(1922):doi.org/10.1038/s41467-021-22202-3.
- 675 [33] Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word represen-
676 tations in vector space. *arXiv*, 1301.3781.
- 677 [34] Moser, K. M., Wei, T., and Brenner, D. (2021). Remote teaching during COVID-19: implications
678 from a national survey of language educators. *System*, 97:102431.
- 679 [35] Nastase, S. A., Goldstein, A., and Hasson, U. (2020). Keep it real: rethinking the primacy of
680 experimental control in cognitive neuroscience. *NeuroImage*, 15(222):117254–117261.

- 681 [36] Nguyen, M., Chang, A., Micciche, E., Meshulam, M., Nastase, S. A., and Hasson, U. (2022).
682 Teacher-student neural coupling during teaching and learning. *Social Cognitive and Affective*
683 *Neuroscience*, 17(4):367–376.
- 684 [37] Poulsen, A. T., Kamronn, S., Dmochowski, J., Parra, L. C., and Hansen, L. K. (2017). EEG
685 in the classroom: synchronised neural recordings during video presentation. *Scientific Reports*,
686 7:43916.
- 687 [38] Ratka, A. (2018). Empathy and the development of affective skills. *American Journal of*
688 *Pharmaceutical Education*, 82(10):doi.org/10.5688/ajpe7192.
- 689 [39] Rhoads, R. A., Berdan, J., and Toven-Lindsey, B. (2013). The open courseware movement in
690 higher education: unmasking power and raising questions about the movement’s democratic
691 potential. *Educational Theory*, 63(1):87–110.
- 692 [40] Scott, P., Asoko, H., and Leach, J. (2007). *Handbook of research on science education*, chapter
693 Student conceptions and conceptual learning in science. Routledge.
- 694 [41] Shao, Y. N., Sun, H. M., Huang, J. W., Li, M. L., Huang, R. R., and Li, N. (2018). Simulation-
695 based empathy training improves the communication skills of neonatal nurses. *Clinical Simula-*
696 *tion in Nursing*, 22:32–42.
- 697 [42] Shim, T. E. and Lee, S. Y. (2020). College students’ experience of emergency remote teaching
698 during COVID-19. *Children and Youth Services Review*, 119:105578.
- 699 [43] Simon, M. A., Tzur, R., Heinz, K., and Kinzel, M. (2004). Explicating a mechanism for
700 conceptual learning: elaborating the construct of reflective abstraction. *Journal for Research in*
701 *Mathematics Education*, 35(5):305–329.
- 702 [44] Stepien, K. A. and Baernstein, A. (2006). Education for empathy. *Journal of General Internal*
703 *Medicine*, 21:524–530.
- 704 [45] Whalen, J. (2020). Should teachers be trained in emergency remote teaching? Lessons learned
705 from the COVID-19 pandemic. *Journal of Technology and Teacher Education*, 28(2):189–199.

- 706 [46] Young, J. R. (2012). Inside the Coursera contract: how an upstart company might profit from
707 free courses. *The Chronicle of Higher Education*, 19(7):1–4.
- 708 [47] Ziman, K., Heusser, A. C., Fitzpatrick, P. C., Field, C. E., and Manning, J. R. (2018). Is
709 automatic speech-to-text transcription ready for use in psychological experiments? *Behavior*
710 *Research Methods*, 50:2597–2605.