

¹ Text embedding models enable high-resolution insights
² into conceptual knowledge and learning in
³ classroom-like settings

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 in real-world educational con-
8 texts. Our approach embeds course content in a high-dimensional conceptual space, where
9 nearby coordinates reflect similar or related concepts. We test our approach using behavioral
10 data from participants who viewed two lecture videos from the Khan Academy platform, inter-
11 leaved between three short multiple-choice quizzes. We applied our framework to the lectures'
12 transcripts and the text of the quiz questions to quantify the conceptual content presented in each
13 moment of video and knowledge probed by each quiz question. We used these embeddings,
14 along with participants' quiz responses, to track how the learners' knowledge changed after
15 watching each video. Our findings demonstrate how a small set of quiz questions may be used
16 to obtain rich and meaningful high-resolution insights into individuals' knowledge, and how it
17 changes over time as they learn.

Introduction

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

30 Designing and building procedures and tools for mapping out knowledge touches on deep
31 questions about what it means to learn. For example, how do we acquire conceptual knowledge?
32 Memorizing course lectures or textbook chapters by rote can lead to the superficial *appearance*
33 of understanding the underlying content, but achieving true conceptual understanding seems
34 to require something deeper and richer. Does conceptual understanding entail connecting newly
35 acquired information to the scaffolding of one’s existing knowledge or experience [6, 10, 13, 14, 56]?
36 Or weaving a lecture’s atomic elements (e.g., its component words) into a structured network
37 that describes how those individual elements are related? Conceptual understanding could also
38 involve building a mental model that transcends the meanings of those individual atomic elements
39 by reflecting the deeper meaning underlying the gestalt whole [33, 37, 53].

40 The difference between “understanding” and “memorizing,” as framed by researchers in ed-
41 ucation, cognitive psychology, and cognitive neuroscience [e.g., 19, 25, 29, 37, 53] has profound
42 analogs in the fields of natural language processing and natural language understanding. For
43 example, considering the raw contents of a document (e.g., its constituent symbols, letters, and

44 words) might provide some information about what the document is about, just as memorizing a
45 passage might provide some ability to answer simple questions about it [e.g., whether it contains
46 words related to furniture versus physics; 7, 8, 36]. However, modern natural language process-
47 ing models [e.g., 9, 11, 44] also attempt to capture the deeper meaning *underlying* those atomic
48 elements. These models consider not only the co-occurrences of those elements within and across
49 documents, but also patterns in how those elements appear across different scales (e.g., sentences,
50 paragraphs, chapters, etc.), the temporal and grammatical properties of the elements, and other
51 high-level characteristics of how they are used [38, 39]. According to these models, the deep
52 conceptual meaning of a document may be captured by a feature vector in a high-dimensional
53 representation space, where nearby vectors reflect conceptually related documents. A model that
54 succeeds at capturing an analog of “understanding” is able to assign nearby feature vectors to
55 two conceptually related documents, *even when the words contained in those documents have very little*
56 *overlap*.

57 Given these insights, what form might the representation of the sum total of a person’s knowl-
58 edge take? First, we might require a means of systematically describing or representing the nearly
59 infinite set of possible things a person could know. Second, we might want to account for potential
60 associations between different concepts. For example, the concepts of “fish” and “water” might be
61 associated in the sense that fish live in water. Third, knowledge may have a critical dependency
62 structure, such that knowing about a particular concept might require first knowing about a set of
63 other concepts. For example, understanding the concept of a fish swimming in water first requires
64 understanding what fish and water *are*. Fourth, as we learn, our “current state of knowledge”
65 should change accordingly. Learning new concepts should both update our characterizations of
66 “what is known” and also unlock any now-satisfied dependencies of those newly learned concepts
67 so that they are “tagged” as available for future learning.

68 Here we develop a framework for modeling how knowledge is acquired during learning. The
69 central idea is to use text embedding models to define the coordinate systems of two maps: (a) a
70 *knowledge map* that describes the extent to which each concept is currently known and (b) a *learning*
71 *map* that describes changes in knowledge over time. Each location on these maps represents

⁷² a single concept, and the maps' geometries are defined such that related concepts are located
⁷³ nearby in space. We use this framework to analyze and interpret behavioral data collected from
⁷⁴ an experiment that had participants watch and answer multiple-choice questions about a series of
⁷⁵ recorded course lectures.

⁷⁶ Our primary research goal is to advance our understanding of what it means to acquire deep,
⁷⁷ real-world conceptual knowledge. Traditional laboratory approaches to studying learning and
⁷⁸ memory (e.g., list learning studies) often draw little distinction between memorization and under-
⁷⁹ standing. Instead, these studies typically focus on whether information is effectively encoded or
⁸⁰ retrieved, rather than whether the information is *understood*. Approaches to studying conceptual
⁸¹ learning, such as category learning experiments, can begin to investigate the distinction between
⁸² memorization and understanding, often by training participants to distinguish arbitrary or ran-
⁸³ dom features in otherwise meaningless categorized stimuli. However the objective of real-world
⁸⁴ training, or learning from life experiences more generally, is often to develop new knowledge that
⁸⁵ may be applied in *useful* ways in the future. In this sense, the gap between modern learning theo-
⁸⁶ ries and modern pedagogical approaches and classroom learning strategies is enormous: most of
⁸⁷ our theories about *how* people learn are inspired by experimental paradigms and models that have
⁸⁸ only peripheral relevance to the kinds of learning that students and teachers actually seek [25, 37].
⁸⁹ To help bridge this gap, our study uses course materials from real online courses to inform, fit,
⁹⁰ and test models of real-world conceptual learning. We also provide a demonstration of how our
⁹¹ models can be used to construct “maps” of what students know, and how their knowledge changes
⁹² with training. In addition to helping to visualize knowledge (and changes in knowledge), we hope
⁹³ that such maps might lead to real-world tools for improving how we educate.

⁹⁴ Results

⁹⁵ At its core, our main modeling approach is based around a simple assumption that we sought to
⁹⁶ test empirically: all else being equal, knowledge about a given concept is predictive of knowledge
⁹⁷ about similar or related concepts. From a geometric perspective, this assumption implies that

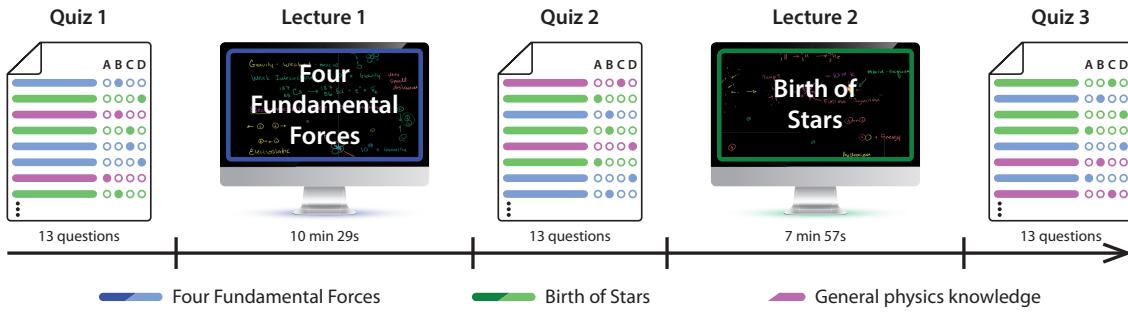


Figure 1: Experimental paradigm. Participants alternate between answering 13-question multiple-choice quizzes and watching two Khan Academy lectures. Each quiz contains a mix of 5 questions about lecture 1, 5 questions about lecture 2, and 3 questions about general physics knowledge. The specific questions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

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 throughout that space. To begin to test this smoothness assumption, we sought to track participants’ knowledge and how it changed over time in response to training.

We asked participants in our study to complete brief multiple-choice quizzes before, between, and after watching two lecture videos from the Khan Academy [32] platform (Fig. 1). The first lecture video, entitled *Four Fundamental Forces*, discussed the four fundamental forces in physics: gravity, strong and weak interactions, and electromagnetism. The second, entitled *Birth of Stars*, provided an overview of our current understanding of how stars form. We selected these lessons to be (a) accessible to a broad audience, i.e., requiring minimal prerequisite knowledge to understand; (b) conceptually related to each other, i.e., covering at least *some* similar or overlapping content; and (c) largely independent of each other, i.e., focused on sufficiently different material that understanding one did not require having seen the other. The two videos we selected are introductory lectures that both belong to Khan Academy’s “Cosmology and Astronomy” course domain, but are taken from different lecture series (“Scale of the Universe” and “Stars, Black Holes, and Galaxies” for the first and second lectures, respectively).

We then created a pool of multiple-choice quiz questions that would enable us to test partici-

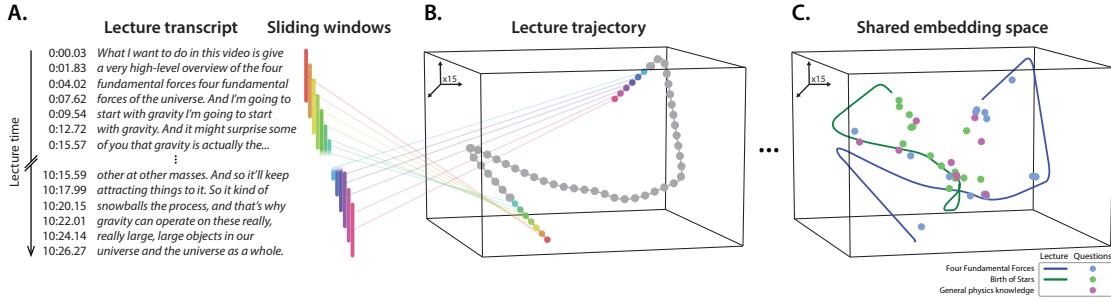


Figure 2: Modeling course content. **A. Building a document pool from sliding windows of text.** We decompose each lecture’s transcript into a series of overlapping sliding windows. The set of transcript snippets (across all windows) may be treated as a set of “documents” for training a text embedding 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.

116 pants’ knowledge about each individual lecture, as well as related content not specifically presented
 117 in either video (see Tab. S1). Participants answered questions randomly drawn from each content
 118 area (lecture 1, lecture 2, and general physics knowledge) on each of the three quizzes. Quiz 1 was
 119 intended to assess participants’ “baseline” knowledge before training, quiz 2 assessed knowledge
 120 after watching the *Four Fundamental Forces* video (i.e., lecture 1), and quiz 3 assessed knowledge
 121 after watching the *Birth of Stars* video (i.e., lecture 2).

122 To study how participants’ conceptual knowledge changed over the course of the experiment,
 123 we first sought to characterize the abstract concepts presented to them in each of the two lectures.
 124 We adapted an approach we developed in prior work [26] to extract the latent themes from the
 125 lectures’ contents using a topic model [8]. Briefly, topic models take as input a collection of text doc-
 126 uments and learn a set of “topics” (i.e., latent themes) from their contents. Once fit, a topic model
 127 can be used to transform arbitrary (potentially new) documents into sets of “topic proportions,” de-
 128 scribing the weighted blend of learned topics reflected in their texts. We parsed automatically gen-
 129 erated transcripts of the two lectures into overlapping sliding windows, which we treated as doc-
 130 uments to fit our model (Fig. 2A; see *Constructing text embeddings of multiple lectures and questions*).

131 Transforming these windows with our model yielded a number-of-windows by number-of-topics
132 (15) topic-proportions matrix, denoting the unique mixture of broad themes from both lectures
133 reflected in each window’s content. Intuitively, each window’s “topic vector” (i.e., column of the
134 topic-proportions matrix) is analogous to a coordinate in a 15-dimensional space (whose axes are
135 topics discovered by the model). Within this space, each lecture’s sequence of topic vectors (i.e.,
136 corresponding to sliding windows parsed from its transcript) forms a *trajectory* that captures how
137 its conceptual content unfolds over time (Fig. 2B). We resampled these trajectories to a resolution
138 of 1 topic vector for each second of video.

139 Next, we sought to characterize what information participants’ performance on each of the
140 three quizzes could provide about their conceptual knowledge at that point in time. Traditional
141 approaches to evaluating students’ performance on short-form knowledge assessments, such as
142 those in our study, entail computing the proportion of correctly answered questions to assign a
143 simple numeric score. While this may afford some measure of the *extent* of a learner’s knowledge,
144 such a score provides little information to either the teacher or the learner about the particular
145 *contents* of their knowledge—in other words, raw “proportion-correct” measures may capture *how*
146 *much* a student knows, but not *what* they know. For instance, suppose a participant in our study was
147 highly knowledgeable about fundamental physical forces (i.e., lecture 1 content) but unfamiliar
148 with how stars are formed (i.e., lecture 2 content), while a second participant was unfamiliar
149 with fundamental forces but had extensive prior knowledge about star formation. Since quiz
150 1 (completed before viewing either lecture) contained an equal number of questions about each
151 lecture’s content (see Fig. 1), these two participants could easily achieve the same proportion-correct
152 score despite their conceptual knowledge differing substantially. How might we distinguish these
153 individuals?

154 We hypothesized that our text embedding model, which we had trained on transcripts of the two
155 lectures to characterize their conceptual contents, should also allow us to capture the conceptual
156 knowledge probed by each quiz question. If our model successfully represented information about
157 the deeper conceptual content of the lectures (i.e., beyond surface-level details such as particular
158 word choices) then we should be able to recover a correspondence between their embeddings and

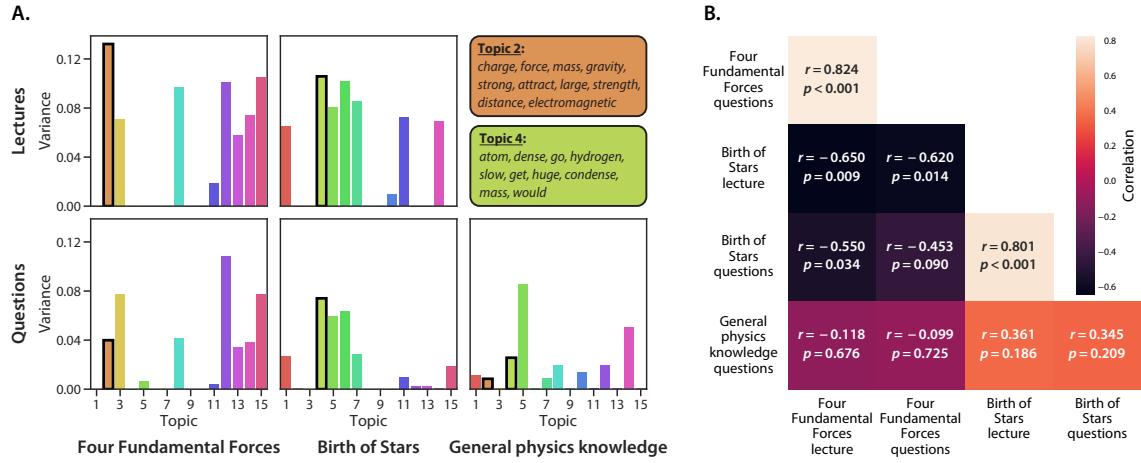


Figure 3: Lecture and question topic overlap. **A. Topic weight variability.** The bar plots display the variance of each topic’s weight across lecture timepoints (top row) and questions (bottom row); colors denote topics. The top-weighted words from the most “expressive” (i.e., variable across observations) 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 category. Each row/column corresponds to a bar plot in panel A.

embeddings of other text that reflects related concepts (e.g., quiz questions that were ostensibly, by design, “about” one of the two lecture) despite differences in exact verbiage. Intuitively, while the content of each given lecture will naturally vary from moment to moment, the particular combination of dimensions along which its trajectory varies across moments correspond to topics that encode meaningful “information” [18] about that lecture’s dynamic conceptual content. We therefore expected that quiz questions pertaining to those same concepts might express the same set of topics.

We used our model to transform the text of each question in our pool into the same embedding space as the two lectures’ trajectories (Fig. 2C). This yielded a single 15-dimensional coordinate (i.e., topic vector) for each question. We then computed the variance of each topic’s weight across timepoints of each lecture, and across questions from each category (i.e., lecture 1-related, lecture 2-related, and general physics knowledge; Fig. 3A). Visual inspection of Figure 3A suggests a strong correspondence between the sets of topics expressed by each lecture and its related questions, with only limited overlap in the topics expressed by non-matched lecture-question

173 set pairs. To quantify these apparent similarities and differences, we computed the correlation
174 between each pair of topic-weight variance distributions (Fig. 3B). We found that our model
175 represented the contents of lecture-related questions using a combination of topic dimensions
176 (and relative variation along those topic dimensions) that was highly similar to that of their
177 reference lecture (*Four Fundamental Forces* (FFF) questions vs. lecture: Pearson's $r(13) = 0.824$, $p <$
178 0.001 , 95% confidence interval (CI) = [0.696, 0.973]; *Birth of Stars* (BoS) questions vs. lecture: $r(13) =$
179 0.801 , $p < 0.001$, 95% CI = [0.539, 0.958]) but diverged significantly from that of the non-reference
180 lecture (FFF questions vs. BoS lecture: $r(13) = -0.620$, $p = 0.014$, 95% CI = [-0.871, -0.326]; BoS
181 questions vs. FFF lecture: $r(13) = -0.550$, $p = 0.034$, 95% CI = [-0.803, -0.246]). This indicated
182 that our model captured the conceptual contents of the lectures and quiz questions sufficiently
183 well to differentiate between questions relating to one lecture versus the other.

184 This could enable us to estimate participants' knowledge for each of the two lectures, separately,
185 by

186 This aspect of our model could enable us to assess participants' knowledge for each lecture
187 separately, based on which *specific* questions they answered right or wrong on each quiz, providing
188 a slightly more refined estimate of the concepts they do and do not know than single proportion-
189 correct score. However, while lectures are often organized around a single, broad theme at a coarse
190 timescale, they typically cover...

191 **Paxton stopped here**

192 Although an individual lecture may be organized around a single broad theme at a coarse
193 scale, at a finer scale each moment of a lecture typically covers a narrower range of content. We
194 wondered whether a text embedding model trained on the lectures' transcripts might capture
195 some of this finer scale content. For example, if a particular question asks about the content
196 from one small part of a lecture, we wondered whether our text embedding model could be used
197 to automatically identify the "matching" moment(s) in the lecture. When we correlated each
198 question's topic vector with the topic vectors for each second of the lectures, we found some
199 evidence that each question is temporally specific (Fig. 4). In particular, most questions' topic
200 vectors were maximally correlated with a well-defined (and relatively narrow) range of timepoints

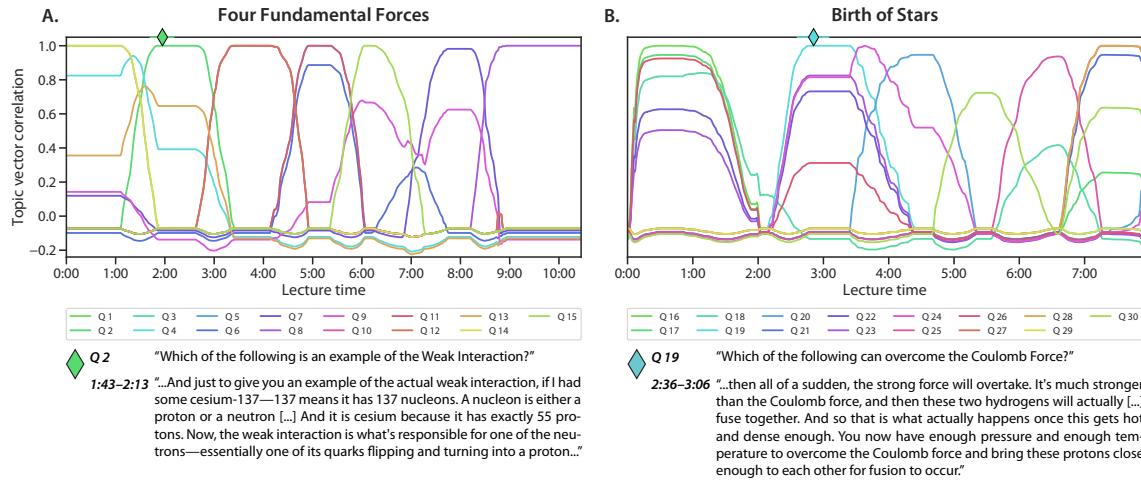


Figure 4: Which parts of each lecture are captured by each question? Each panel displays timeseries plots showing how each question’s topic vector correlates with each video timepoint’s topic vector (Panel A.: correlations for the *Four Fundamental Forces* lecture and associated questions; Panel B.: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

from their corresponding lectures, and the correlations fell off sharply outside of that range. We also examined the best-matching intervals for each question qualitatively by comparing the text of the question to the text of the most-correlated parts of the lectures. Despite that the questions were excluded from the text embedding model’s training set, in general we found (through manual inspection) a close correspondence between the conceptual content that each question covered and the content covered by the best-matching moments of the lectures. Two representative examples are shown at the bottom of Figure 4.

The ability to quantify how much each question is “asking about” the content from each moment of the lectures could enable high-resolution insights into participants’ knowledge. Traditional approaches to estimating how much a student “knows” about the content of a given lecture entail computing the proportion of correctly answered questions. But if two students receive identical scores on an exam, might our 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 fill in the “gaps” in the two students’ understandings, we might do well to focus on concept *A* for the first student, but to also add in materials pertaining to concepts *B* and *C* for the second student.

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 the coordinate is given by the weighted average 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 timecourse describing how much “knowledge” the participant has about any part of the lecture. As shown in Figure 5, we can also apply this approach separately for the questions from each quiz the

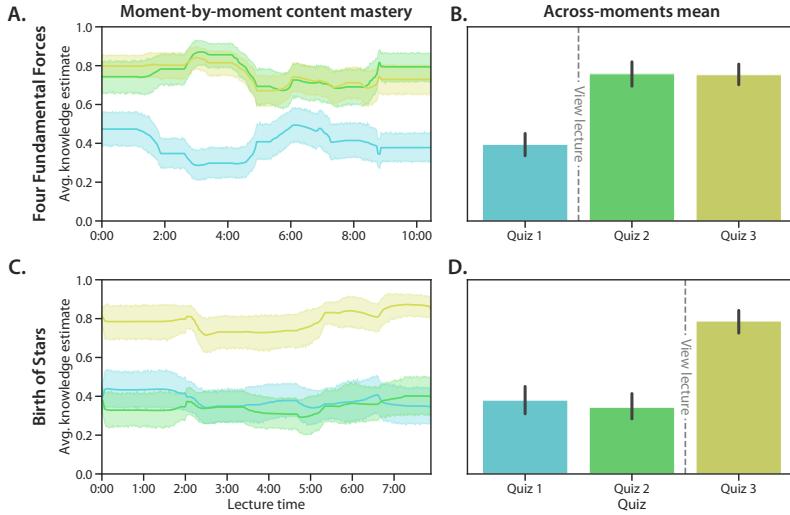


Figure 5: Estimating moment-by-moment knowledge acquisition. **A. Moment-by-moment knowledge about the *Four Fundamental Forces*.** Each trace displays the weighted proportion of correctly answered questions about the content reflected in each moment of the lecture (see *Estimating dynamic knowledge traces*), using responses from one quiz’s 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.

229 participants took throughout the experiment. From just 13 questions per quiz, we obtain a high-
 230 resolution snapshot (at the time each quiz was taken) of what the participants knew about any
 231 moment’s content, from either of the two lectures they watched (comprising a total of 1106 samples
 232 across the two lectures).

233 Of course, even though the timecourses in Figure 5A and C provide detailed *estimates* about
 234 participants’ knowlege, those estimates are only *useful* to the extent that they accurately reflect what
 235 participants actually know. As one sanity check, we anticipated that the knowledge estimates
 236 should show a content-specific “boost” in participants’ knowledge after watching each lecture.
 237 In other words, if participants learn about each lecture’s content when they watch each lecture,
 238 the knowledge estimates should reflect that. After watching the *Four Fundamental Forces* lecture,
 239 participants should show 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 the *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 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, the estimated knowledge should have some predictive information about whether the participant is likely to answer the question correctly or incorrectly. 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



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.

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 in estimated knowledge for held-out correctly answered questions was larger than for held-out incorrectly answered 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

284 text embedding space.

285 We projected the resulting 100-dimensional topic vectors (for each second of video and for each
286 question) into a shared 2-dimensional space (see *Creating knowledge and learning map visualizations*).
287 Next, we sampled points evenly from a 100×100 grid of coordinates that evenly tiled a rectangle
288 enclosing the 2D projections of the videos and questions. We used Equation 4 to estimate partici-
289 pants' knowledge at each of these 10,000 sampled locations, and we averaged these estimates across
290 participants to obtain an estimated average *knowledge map* (Fig. 7). Intuitively, the knowledge map
291 constructed from a given quiz's responses provides a visualization of how "much" participants
292 know about any content expressible by the fitted text embedding model.

293 Several features of the resulting knowledge maps are worth noting. The average knowledge
294 map estimated from Quiz 1 responses (Fig. 7, leftmost map) shows that participants tended to
295 have relatively little knowledge about any parts of the text embedding space (i.e., the shading
296 is relatively dark everywhere). The knowledge map estimated from Quiz 2 responses shows a
297 marked increase in knowledge on the left side of the map (around roughly the same range of
298 coordinates covered by the *Four Fundamental Forces* lecture, indicated by the dotted blue line).
299 In other words, participants' estimated increase in knowledge is localized to conceptual content
300 that is nearby (i.e., related to) the content from the lecture they watched prior to taking Quiz
301 2. This localization is non-trivial: the knowledge estimates are informed only by the embedded
302 coordinates of the *quiz questions*, not by the embeddings of either lecture (Eqn. 4). Finally, the
303 knowledge map estimated from Quiz 3 responses shows a second increase in knowledge, localized
304 to the region surrounding the embedding of the *Birth of Stars* lecture participants watched prior to
305 taking Quiz 3.

306 Another way of visualizing these content-specific increases in knowledge (apparently driven
307 by watching each lecture) is displayed in Figure 7B. Taking the point-by-point difference between
308 the knowledge maps estimated from responses to a successive pair of quizzes yields a *learning map*
309 that describes the *change* in knowledge estimates from one quiz to the next. These learning maps
310 highlight that the estimated knowledge increases we observed across maps were specific to the
regions around the embeddings of each lecture in turn.



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 and the coordinates of each question are indicated by dotted lines and dots. Each map reflects an average across all participants. For individual participants’ maps, see Figures S1, S2, and S3. **B.** Average “learning maps” estimated between each successive pair of quizzes. The learning maps are in the same general format as the knowledge maps in Panel A, but each coordinate in the learning maps indicates the difference between the corresponding coordinates in the indicated pair of knowledge maps—i.e., how much the estimated knowledge “changed” across the two quizzes. Each map reflects an average across all participants. For individual participants’ maps, see Figures S4 and S5. **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 in 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.

Because the 2D projection we used to construct the knowledge and learning maps is invertible, we may gain additional insights into the estimates by reconstructing the original high-dimensional topic vectors for any point(s) in the maps we are interested in. For example, this could serve as a useful tool for an instructor looking to better understand which content areas a student (or a group of students) knows well (or poorly). As a demonstration, we show the top-weighted words from the blends of topics reconstructed from three example locations on the maps (Fig. 7C): one point near the *Four Fundamental Forces* embedding (yellow); a second point near the *Birth of Stars* embedding (orange), and a third point somewhere in between the two lectures' embeddings (pink). As shown in the word clouds in the Panel, the top-weighted words at the example coordinate near the *Four Fundamental Forces* embedding also tended to be weighted heavily by the topics expressed in that lecture. Similarly, the top-weighted words at the example coordinate near the *Birth of Stars* embedding tended to be weighted most heavily by the topics expressed in *that* lecture. And the top-weighted words at the example coordinate between the two lectures' embeddings show a roughly even mix of words most strongly associated with each lecture.

Discussion

Teaching, like effective writing and speaking, is fundamentally about empathy [1, 43, 58]. Great teachers consider students' interests [12, 59], backgrounds [15, 46, 52], and working memory capacities [2], and flexibly optimize their teaching strategies within those constraints [4, 22, 27]. In the classroom, empathizing with students also means maintaining open lines of communication [64] by fostering an environment in which all students feel comfortable speaking up if they have an exciting new idea, or if they are having trouble understanding something [20, 60]. In-person instruction also often entails dynamic student-teacher and student-student interactions. These in-person interactions can provide the instructor with valuable information about students' understanding of the course material, beyond what they can glean solely from exams or assignments [17, 24, 61]. In turn, this can allow the instructor to adapt their teaching approaches on-the-fly according to students' questions and behaviors. But what does great teaching look like in asynchronous online

338 courses, when the instructor typically prepares course lectures and materials without knowing
339 who will ultimately be learning from them? Can the empathetic side of teaching be automated
340 and scaled?

341 The notion of empathy also related to “theory of mind” of other individuals [21, 28, 41].
342 Considering others’ unique perspectives, prior experiences, knowledge, goals, etc., can help us
343 to more effectively interact and communicate [50, 54, 57]. The knowledge and learning maps
344 we estimate in our study (Fig. 7) hint at one potential form that an automated “empathetic”
345 teacher might take. We imagine automated content delivery systems that adapt lessons on the
346 fly according to continually updated estimates of what students know and how quickly they are
347 learning different conceptual content [e.g., building on ideas such as 3, 23, 35, 63, and others].

348 Over the past several years, the global pandemic has forced many educators to teach re-
349 motely [31, 45, 55, 62]. This change in world circumstances is happening alongside (and perhaps
350 accelerating) geometric growth in the availability of high quality online courses on platforms such
351 as Khan Academy [32], Coursera [65], EdX [34], and others [51]. Continued expansion of the global
352 internet backbone and improvements in computing hardware have also facilitated improvements
353 in video streaming, enabling videos to be easily downloaded and shared by large segments of the
354 world’s population. This exciting time for online course instruction provides an opportunity to
355 re-evaluate how we, as a global community, educate ourselves and each other. For example, we
356 can ask: what makes an effective course or training program? Which aspects of teaching might be
357 optimized or automated? How can we provide How and why do learning needs and goals vary
358 across people? How might we lower barriers to achieving a high quality education?

359 Alongside these questions, there is a growing desire to extend existing theories beyond the
360 domain of lab testing rooms and into real classrooms [30]. In part, this has led to a recent
361 resurgence of “naturalistic” or “observational” experimental paradigms that attempt to better
362 reflect more ethologically valid phenomena that are more directly relevant to real-world situations
363 and behaviors [47]. In turn, this has brought new challenges in data analysis and interpretation. A
364 key step towards solving these challenges will be to build explicit models of real-world scenarios
365 and how people behave in them (e.g., models of how people learn conceptual content from real-

366 world courses, as in our current study). A second key step will be to understand which sorts
367 of signals derived from behaviors and/or other measurements [e.g., neurophysiological data; 5,
368 16, 42, 48, 49] might help to inform these models. A third major step will be to develop and
369 employ reliable ways of evaluating the complex models and data that are a hallmark of naturalistic
370 paradigms.

371 Ultimately, our work suggests a new line of questions regarding the future of education:
372 which aspects of teaching can be optimized and/or automated? The social benefits of face-to-face
373 instruction, such as social interactions, friendships, and emotional support, cannot (and perhaps
374 should not) be fully replaced by an automated computer-based system. Nor can modern computer
375 systems experience emotional empathy in the human sense of the word. On the other hand,
376 perhaps it is possible to separate out the social aspects of classroom instruction from the purely
377 learning-related aspects. Our study shows that text embedding models can uncover detailed
378 insights into students' knowledge and how it changes over time during learning. We hope that
379 these advances might help pave the way for new ways of teaching or delivering educational content
380 that are tailored to individual students' learning needs and goals.

381 Materials and methods

382 Participants

383 We enrolled a total of 50 Dartmouth undergraduate students in our study. Participants received
384 course credit for enrolling. We asked each participant to fill out a demographic survey that included
385 questions about their age, gender, native spoken language, ethnicity, race, hearing, color vision,
386 sleep, coffee consumption, level of alertness, and several aspects of their educational background
387 and prior coursework.

388 Participants' ages ranged from 18 to 22 years (mean: 19.52 years; standard deviation: 1.09
389 years). A total of 15 participants reported their gender as male and 35 participants reported their
390 gender as female. A total of 49 participants reported their native language as "English" and 1

391 reported having another native language. A total of 47 participants reported their ethnicity as
392 “Not Hispanic or Latino” and three reported their ethnicity as “Hispanic or Latino.” Participants
393 reported their races as White (32 participants), Asian (14 participants), Black or African American
394 (5 participants), American Indian or Alaska Native (1 participant), and Native Hawaiian or Other
395 Pacific Islander (1 participant). (Note that some participants selected multiple racial categories.)

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

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

408 Participants reported their undergraduate major(s) as “social sciences” (28 participants), “nat-
409 ural sciences” (16 participants), “professional” (e.g., pre-med or pre-law; 8 participants), “mathe-
410 matics and engineering” (7 participants), “humanities” (4 participants), or “undecided” (3 partici-
411 pants). Note that some participants selected multiple categories for their undergraduate major. We
412 also asked participants about the courses they had taken. In total, 45 participants reported having
413 taken at least one Khan Academy course in the past, and 5 reported not having taken any Khan
414 Academy courses. Of those who reported having watched at least one Khan Academy course,
415 7 participants reported having watched 1–2 courses, 11 reported having watched 3–5 courses, 8
416 reported having watched 5–10 courses, and 19 reported having watched 10 or more courses. We
417 also asked participants about the specific courses they had watched, categorized under different
418 subject areas. In the “Mathematics” area, participants reported having watched videos on AP

419 Calculus AB (21 participants), Precalculus (17 participants), Algebra 2 (14 participants), AP Cal-
420 culus BC (12 participants), Trigonometry (11 participants), Algebra 1 (10 participants), Geometry
421 (8 participants), Pre-algebra (7 participants), Multivariable Calculus (5 participants), Differential
422 Equations (5 participants), Statistics and Probability (4 participants), AP Statistics (2 participants),
423 Linear Algebra (2 participants), Early Math (1 participant), Arithmetic (1 participant), and other
424 videos not listed in our survey (5 participants). In the “Science and engineering” area, participants
425 reported having watched videos on Chemistry, AP Chemistry, or Organic Chemistry (21 partic-
426 ipants); Physics, AP Physics I, or AP Physics II (15 participants); Biology, AP Biology; or High
427 school Biology (15 participants); Health and Medicine (1 participant); or other videos not listed
428 in our survey (19 participants). We also asked participants whether they had specifically seen the
429 videos used in our experiment. Of the 45 participants who reported having taken at least
430 one Khan Academy course in the past, 44 participants reported that they had not watched the *Four*
431 *Fundamental Forces* video, and 1 participant reported that they were not sure whether they had
432 watched it. All participants reported that they had not watched the *Birth of Stars* video. When
433 we asked participants about non-Khan Academy online courses, they reported having watched
434 or taken courses on Mathematics (15 participants), Science and engineering (11 participants), Test
435 preparation (9 participants), Economics and finance (3 participants), Arts and humanities (2 partic-
436 ipants), Computing (2 participants), and other categories not listed in our survey (18 participants).
437 Finally, we asked participants about in-person courses they had taken in different subject areas.
438 They reported taking courses in Mathematics (39 participants), Science and engineering (38 par-
439 ticipants), Arts and humanities (35 participants), Test preparation (27 participants), Economics
440 and finance (26 participants), Computing (15 participants), College and careers (7 participants), or
441 other courses not listed in our survey (6 participants).

442 **Experiment**

443 We hand-selected two course videos from the Khan Academy platform: *Four Fundamental Forces*
444 (an introduction to gravity, electromagnetism, the weak nuclear force, and the strong nuclear force;
445 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 after viewing lecture 2 (Fig. 1). The questions appearing on each quiz, for each participant, were randomly chosen from the full set of 39, with the constraints that (a) each quiz contain 5 questions about lecture 1, 5 questions about lecture 2, and 3 questions about general physics knowledge, and (b) each question appear exactly once for each participant. The order of questions on each quiz, and the order of answer options for each question, were also randomized. Our experimental protocol was approved by the Committee for the Protection of Human Subjects at Dartmouth College. We used the experiment to develop and test our computational framework for estimating knowledge and learning.

Analysis

Constructing text embeddings of multiple lectures and questions

We adapted an approach we developed in prior work [26] to embed each moment of the two lectures and each question in our pool in a common representational space. Briefly, our approach uses a topic model (Latent Dirichlet Allocation; 8), trained on a set of documents, to discover a set of k “topics” or “themes.” Formally, each topic is defined as a set of weights over each word in the model’s vocabulary (i.e., the union of all unique words, across all documents, excluding “stop words.”). Conceptually, each topic is intended to give larger weights to words that are conceptually related or that tend to co-occur in the same documents. After fitting a topic model, each document in the training set, or any *new* document that contains at least some of the words in the model’s vocabulary, may be represented as a k -dimensional vector describing how much the

472 document (most probably) reflects each topic. (Unless, otherwise noted, we used $k = 15$ topics.)

473 As illustrated in Figure 2A, we start by building up a corpus of documents using overlapping
474 sliding windows that span each video’s transcript. Khan Academy provides professionally created,
475 manual transcriptions of all videos for closed captioning. However, such transcripts would not
476 be readily available in all contexts to which our framework could potentially be applied. Khan
477 Academy videos are hosted on the YouTube platform, which additionally provides automated
478 captions. We opted to use these automated transcripts (which, in prior work, we have found are
479 sufficiently near-human quality yield reliable data in behavioral studies; 66) when developing our
480 framework in order to make it more easily extensible and adaptable by others in the future.

481 We fetched these automated transcripts using the youtube-transcript-api Python package
482 (**Jeremy can you add a citation for this?** <https://github.com/jdepoix/youtube-transcript-api>). The
483 transcripts consisted of one timestamped line of text for every few seconds (mean: 2.34s; standard
484 deviation: 0.83s) of spoken content in the video (i.e., corresponding to each individual caption that
485 would appear on-screen if viewing the lecture via YouTube, and when those lines would appear).
486 We defined a sliding window length of (up to) $w = 30$ transcript lines, and assigned each window
487 a timestamp corresponding to the midpoint between its first and last lines’ timestamps. These
488 sliding windows ramped up and down in length at the very beginning and end of the transcript,
489 respectively. In other words, the first sliding window covered only the first line from the transcript;
490 the second sliding window covered the first two lines; and so on. This insured that each line of
491 the transcript appeared in the same number (w) of sliding windows. After performing various
492 standard text preprocessing (e.g., normalizing case, lemmatizing, removing punctuation and stop-
493 words), we treated the text from each sliding window as a single “document,” and we combined
494 these documents across the two videos’ windows to create a single training corpus for the topic
495 model. The top words from each of the 15 discovered topics may be found in Table S2.

496 After fitting a topic model to each videos’ transcripts, we could use the trained model to
497 transform arbitrary (potentially new) documents into k -dimensional topic vectors. A convenient
498 property of these topic vectors is that documents that reflect similar blends of topics (i.e., documents
499 that reflect similar themes, according to the model) will yield similar (in terms of Euclidean **should**

500 we say something like cosine instead of euclidean? distance, correlation, etc.) topic vectors. In
501 general, the similarity between different documents' topic vectors may be used to characterize the
502 similarity in conceptual content between the documents.

503 We transformed each sliding window's text into a topic vector, and then used linear interpolation
504 (independently for each topic dimension) to resample the resulting timeseries to one vector
505 per second. We also used the fitted model to obtain topic vectors for each question in our pool
506 (Tab. S1). Taken together, we obtained a *trajectory* for each video, describing its path through topic
507 space, and a single coordinate for each question (Fig. 2C). Embedding both videos and all of the
508 questions using a common model enables us to compare the content from different moments of
509 videos, compare the content across videos, and estimate potential associations between specific
510 questions and specific moments of video.

511 **Estimating dynamic knowledge traces**

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

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

514 where

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

515 and where mincorr and maxcorr are the minimum and maximum correlations between any lecture
516 timepoint and question, taken over all timepoints and questions across both lectures and all five
517 question used to estimate the knowledge trace **Note: make sure this is correct in results section;**
518 **not sure i caught all instances**. We also define $f(s, \Omega)$ as the s^{th} topic vector from the set of topic
519 vectors Ω . Here t indexes the set of lecture topic vectors, L , and i and j index the topic vectors of
520 questions used to estimate the knowledge trace, Q . Note that "correct" denotes the set of indices
521 of the questions the participant answered correctly on the given quiz.

522 Intuitively, $\text{ncorr}(x, y)$ is the correlation between two topic vectors (e.g., the topic vector from one

523 timepoint in a lecture, x , and the topic vector for one question, y), normalized by the minimum and
524 maximum correlations (across all timepoints t and questions Q) to range between 0 and 1, inclusive.
525 Equation 1 then computes the weighted average proportion of correctly answered questions about
526 the content presented at timepoint t , where the weights are given by the normalized correlations
527 between timepoint t 's topic vector and the topic vectors for each question. The normalization
528 step (i.e., using `ncorr` instead of the raw correlations) insures that every question contributes some
529 non-zero amount to the knowledge estimate.

530 **Creating knowledge and learning map visualizations**

531 An important feature of our approach is that, given a trained text embedding model and partic-
532 ipants' quiz performance on each question, we can estimate their knowledge about *any* content
533 expressible by the embedding model—not solely the content explicitly probed by the quiz ques-
534 tions or even appearing in the lectures. To visualize these estimates (Figs. 7, S1, S2, S3, S4, and S5),
535 we used Uniform Manifold Approximation and Projection (UMAP; 40) to construct a 2D projection
536 of the text embedding space. Sampling the original 100-dimensional space at high resolution to
537 obtain an adequate set of topic vectors spanning the embedding space would be computationally
538 intractable. However, sampling a 2D grid is much more feasible. We defined a rectangle enclosing
539 the 2D projections of the lectures' and quizzes' embeddings, and we sampled points from a regular
540 100×100 grid of coordinates that evenly tiled the enclosing rectangle. We sought to estimate
541 participants' knowledge (and learning, i.e., changes in knowledge) at each of the resulting 10,000
542 coordinates.

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

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

547 The λ term in the RBF equation controls the “smoothness” of the function, where larger values

548 of λ result in smoother maps. In our implementation we used $\lambda = 50$. Next, we estimated the
549 “knowledge” at each coordinate, x , using:

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

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

554 **Author contributions**

555 Conceptualization: PCF, ACH, and JRM. Methodology: PCF, ACH, and JRM. Software: PCF.
556 Validation: PCF. Formal analysis: PCF. Resources: PCF, ACH, and JRM. Data curation: PCF.
557 Writing (original draft): JRM. Writing (review and editing): PCF, ACH, and JRM. Visualization:
558 PCF and JRM. Supervision: JRM. Project administration: PCF. Funding acquisition: JRM.

559 **Data and code availability**

560 All of the data analyzed in this manuscript, along with all of the code for running our experiment
561 and carrying out the analyses may be found at <https://github.com/ContextLab/efficient-learning-khan>.

563 **Acknowledgements**

564 We acknowledge useful discussions, assistance in setting up an earlier (unpublished) version of
565 this study, and assistance with some of the data collection efforts from Will Baxley, Max Bluestone,
566 Daniel Carstensen, Kunal Jha, Caroline Lee, Lucy Owen, Xinming Xu, and Kirsten Ziman. Our
567 work was supported in part by NSF CAREER Award Number 2145172 to JRM. The content is

568 solely the responsibility of the authors and does not necessarily represent the official views of our
569 supporting organizations. The funders had no role in study design, data collection and analysis,
570 decision to publish, or preparation of the manuscript.

571 **References**

- 572 [1] Aldrup, K., Carstensen, B., and Klusmann, U. (2022). Is empathy the key to effective teaching?
573 A systematic review of its association with teacher-student interactions and student outcomes.
574 *Educational Psychology Review*, 34:1177001216.
- 575 [2] Alloway, T. P. (2012). Teachers' perceptions of classroom behaviour and working memory.
576 *Educational Research and Review*, 7(6):138–142.
- 577 [3] Anderson, J. R. and Skwarecki, E. (1986). The automated tutoring of introductory computer
578 programming. *Communications of the ACM*, 29(9):842–849.
- 579 [4] Anderton, R. S., Vitali, J., Blackmore, C., and Bakeberg, M. C. (2021). Flexible teaching and learn-
580 ing modalities in undergraduate science amid the COVID-19 pandemic. *Frontiers in Education*,
581 5:doi.org/10.3389/feduc.2020.609703.
- 582 [5] Bevilacque, D., Davidesco, I., Wan, L., and Chaloner, K. (2019). Brain-to-brain synchrony and
583 learning outcomes vary by student-teacher dynamics: evidence from a real-world classroom
584 electroencephalography study. *Journal of Cognitive Neuroscience*, 31(3):401–411.
- 585 [6] Blaye, A., Bernard-Peyron, V., Paour, J.-L., and Bonthoux, F. (2006). Category flexibility in chil-
586 dren: distinguishing response flexibility from conceptual flexibility; the protracted development
587 of taxonomic representations. *European Journal of Developmental Psychology*, 3(2):163–188.
- 588 [7] Blei, D. M. and Lafferty, J. D. (2006). Dynamic topic models. In *Proceedings of the International*
589 *Conference on Machine Learning*, pages 113–120, New York, NY. Association for Computing
590 Machinery.

- 591 [8] Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine*
592 *Learning Research*, 3:993–1022.
- 593 [9] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A.,
594 Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child,
595 R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M.,
596 Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei,
597 D. (2020). Language models are few-shot learners. *arXiv*, 2005.14165.
- 598 [10] Caramazza, A. and Mahon, B. Z. (2003). The organization of conceptual knowledge: the
599 evidence from category-specific semantic deficits. *Trends in Cognitive Sciences*, 7(8):354–361.
- 600 [11] Cer, D., Yang, Y., Kong, S. Y., Hua, N., Limtiaco, N., John, R. S., Constant, N., Guajardo-
601 Cespedes, M., Yuan, S., Tar, C., Sung, Y.-H., Strope, B., and Kurzweil, R. (2018). Universal
602 sentence encoder. *arXiv*, 1803.11175.
- 603 [12] Clark, J. (2010). Powerpoint and pedagogy: maintaining student interest in university lectures.
604 *College Teaching*, 56(1):39–44.
- 605 [13] Constantinescu, A. O., O'Reilly, J. X., and Behrens, T. E. J. (2016). Organizing conceptual
606 knowledge in humans with a gridlike code. *Science*, 352(6292):1464–1468.
- 607 [14] Deacon, D., Grose-Fifer, J., Yang, C. M., Stanick, V., Hewitt, S., and Dynowska, A. (2004).
608 Evidence for a new conceptualization of semantic representation in the left and right cerebral
609 hemispheres. *Cortex*, 40(3):467–478.
- 610 [15] den Brok, P., van Tartwijk, J., Wubbels, T., and Veldman, I. (2010). The differential effect of
611 the teacher-student interpersonal relationship on student outcomes for students with different
612 ethnic backgrounds. *British Journal of Educational Psychology*, 80(2):199–221.
- 613 [16] Dikker, S., Wan, L., Davidesco, I., Kaggen, L., Oostrik, M., McClintock, J., Rowland, J.,
614 Michalareas, G., van Bavel, J. J., Ding, M., and Poeppel, D. (2017). Brain-to-brain synchrony
615 tracks real-world dynamic group interactions in the classroom. *Current Biology*, 27(9):1375–1380.

- 616 [17] Englehart, J. M. (2009). Teacher-student interaction. In Saha, L. J. and Dworkin, A. G., editors,
617 *International Handbook of Research on Teachers and Teaching*. Springer International Handbooks of
618 Education.
- 619 [18] Fisher, R. A. (1922). On the mathematical foundations of theoretical statistics. *Philosophical
620 Transactions of the Royal Society A*, 222(602):309–368.
- 621 [19] Gallagher, J. J. (2000). Teaching for understanding and application of science knowledge.
622 *School Science and Mathematics*, 100(6):310–318.
- 623 [20] Garran, A. M. and Rasmussen, B. M. (2014). Safety in the classroom: reconsidered. *Journal of
624 Teaching in Social Work*, 34(4):401–412.
- 625 [21] Goldstein, T. R. and Winner, E. (2012). Enhancing empathy and theory of mind. *Journal of
626 Cognition and Development*, 13(1):19–37.
- 627 [22] Goode, S., Willis, R. A., Wolf, J. R., and Harris, A. L. (2007). Enhancing IS education with
628 flexible teaching and learning. *Journal of Information Systems Education*, 18(3):297–302.
- 629 [23] Halff, H. M. (1988). Curriculum and instruction in automated tutors. *Foundations of intelligent
630 tutoring systems*, pages 79–108.
- 631 [24] Hall, J. K. and Walsh, M. (2002). Teacher-student interaction and language learning. *Annual
632 Review of Applied Linguistics*, 22:186–203.
- 633 [25] Hall, R. and Greeno, J. (2008). *21st century education: A reference handbook*, chapter Conceptual
634 learning, pages 212–221. Sage Publications.
- 635 [26] Heusser, A. C., Fitzpatrick, P. C., and Manning, J. R. (2021). Geometric models reveal be-
636 havioral and neural signatures of transforming naturalistic experiences into episodic memories.
637 *Nature Human Behavior*, 5:905–919.
- 638 [27] Johnston, S. (2002). Introducing and supporting change towards more flexible teaching ap-
639 proaches. In *The convergence of distance and conventional education*. Routledge.

- 640 [28] Kanske, P., Böckler, A., and Singer, T. (2015). Models, mechanisms and moderators dissociating
641 empathy and theory of mind. In *Social Behavior from Rodents to Humans*, pages 193–206. Springer.
- 642 [29] Katona, G. (1940). *Organizing and memorizing: studies in the psychology of learning and teaching*.
643 Columbia University Press.
- 644 [30] Kaufman, D. M. (2003). Applying educational theory in practice. *British Medical Journal*,
645 326(7382):213–216.
- 646 [31] Kawasaki, H., Yamasaki, S., Masuoka, Y., Iwasa, M., Fukita, S., and Matsuyama, R. (2021).
647 Remote teaching due to COVID-19: an exploration of its effectiveness and issues. *International
648 Journal of Environmental Research and Public Health*, 18(5):2672.
- 649 [32] Khan, S. (2004). *The Khan Academy*. Salman Khan.
- 650 [33] Kintsch (1970). *Learning, memory, and conceptual processes*. Wiley.
- 651 [34] Kolowich, S. (2013). How EdX plans to earn, and share, revenue from its free online courses.
652 *The Chronicle of Higher Education*, 21:1–5.
- 653 [35] Kumar, A. N. (2005). Generation of problems, answers, grade, and feedback—case study of a
654 fully automated tutor. *Journal on Educational Resources in Computing*, 5(3):1–25.
- 655 [36] Landauer, T. K. and Dumais, S. T. (1997). A solution to Plato’s problem: the latent semantic
656 analysis theory of acquisition, induction, and representation of knowledge. *Psychological Review*,
657 104:211–240.
- 658 [37] Maclellan, E. (2005). Conceptual learning: the priority for higher education. *British Journal of
659 Educational Studies*, 53(2):129–147.
- 660 [38] Manning, J. R. (2020). Context reinstatement. In Kahana, M. J. and Wagner, A. D., editors,
661 *Handbook of Human Memory*. Oxford University Press.
- 662 [39] Manning, J. R. (2021). Episodic memory: mental time travel or a quantum “memory wave”
663 function? *Psychological Review*, 128(4):711–725.

- 664 [40] McInnes, L., Healy, J., and Melville, J. (2018). UMAP: Uniform manifold approximation and
665 projection for dimension reduction. *arXiv*, 1802(03426).
- 666 [41] Meltzoff, A. N. (2011). Social cognition and the origins of imitation, empathy, and theory of
667 mind. In *The Wiley-Blackwell handbook of childhood cognitive development*. Wiley-Blackwell.
- 668 [42] Meshulam, M., Hasenfratz, L., Hillman, H., Liu, Y. F., Nguyen, M., Norman, K. A., and Hasson,
669 U. (2020). Neural alignment predicts learning outcomes in students taking an introduction to
670 computer science course. *Nature Communications*, 12(1922):doi.org/10.1038/s41467-021-22202-3.
- 671 [43] Meyers, S., Rowell, K., Wells, M., and Smith, B. C. (2019). Teacher empathy: a model of
672 empathy for teaching for student success. *College Teaching*, 67(3):160–168.
- 673 [44] Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word represen-
674 tations in vector space. *arXiv*, 1301.3781.
- 675 [45] Moser, K. M., Wei, T., and Brenner, D. (2021). Remote teaching during COVID-19: implications
676 from a national survey of language educators. *System*, 97:102431.
- 677 [46] Muijs, D. and Reynolds, D. (2003). Student background and teacher effects on achievement and
678 attainment in mathematics: a longitudinal study. *Educational Research and Evaluation*, 9(3):289–
679 314.
- 680 [47] Nastase, S. A., Goldstein, A., and Hasson, U. (2020). Keep it real: rethinking the primacy of
681 experimental control in cognitive neuroscience. *NeuroImage*, 15(222):117254–117261.
- 682 [48] Nguyen, M., Chang, A., Micciche, E., Meshulam, M., Nastase, S. A., and Hasson, U. (2022).
683 Teacher-student neural coupling during teaching and learning. *Social Cognitive and Affective
684 Neuroscience*, 17(4):367–376.
- 685 [49] Poulsen, A. T., Kamronn, S., Dmochowski, J., Parra, L. C., and Hansen, L. K. (2017). EEG
686 in the classroom: synchronised neural recordings during video presentation. *Scientific Reports*,
687 7:43916.

- 688 [50] Ratka, A. (2018). Empathy and the development of affective skills. *American Journal of*
689 *Pharmaceutical Education*, 82(10):doi.org/10.5688/ajpe7192.
- 690 [51] Rhoads, R. A., Berdan, J., and Toven-Lindsey, B. (2013). The open courseware movement in
691 higher education: unmasking power and raising questions about the movement's democratic
692 potential. *Educational Theory*, 63(1):87–110.
- 693 [52] Rosenshine, B. (1976). Recent research on teaching behaviors and student achievement. *Journal*
694 *of Teacher Education*, 27(1):61–64.
- 695 [53] Scott, P., Asoko, H., and Leach, J. (2007). *Handbook of research on science education*, chapter
696 Student conceptions and conceptual learning in science. Routledge.
- 697 [54] Shao, Y. N., Sun, H. M., Huang, J. W., Li, M. L., Huang, R. R., and Li, N. (2018). Simulation-
698 based empathy training improves the communication skills of neonatal nurses. *Clinical Simula-*
699 *tion in Nursing*, 22:32–42.
- 700 [55] Shim, T. E. and Lee, S. Y. (2020). College students' experience of emergency remote teaching
701 during COVID-19. *Children and Youth Services Review*, 119:105578.
- 702 [56] Simon, M. A., Tzur, R., Heinz, K., and Kinzel, M. (2004). Explicating a mechanism for
703 conceptual learning: elaborating the construct of reflective abstraction. *Journal for Research in*
704 *Mathematics Education*, 35(5):305–329.
- 705 [57] Stepien, K. A. and Baernstein, A. (2006). Education for empathy. *Journal of General Internal*
706 *Medicine*, 21:524–530.
- 707 [58] Stojiljković, S., Djigić, G., and Zlatković, B. (2012). Empathy and teachers' roles. *Procedia –*
708 *Social and Behavioral Sciences*, 69:960–966.
- 709 [59] Swarat, S., Ortony, A., and Revelle, W. (2012). Activity matters: understanding student interest
710 in school science. *Journal of Research in Science Teaching*, 49(4):515–537.
- 711 [60] Turner, S. and Braine, M. (2015). Unravelling the 'safe' concept in teaching: what can we learn
712 from teachers' understanding? *Pastoral Care in Education*, 33(1):47–62.

- 713 [61] van de Pol, J., Volman, M., and Beishuizen, J. (2010). Scaffolding in teacher-student interaction:
714 a decade of research. *Educational Psychology Review*, 22:271–296.
- 715 [62] Whalen, J. (2020). Should teachers be trained in emergency remote teaching? Lessons learned
716 from the COVID-19 pandemic. *Journal of Technology and Teacher Education*, 28(2):189–199.
- 717 [63] Wolz, U., McKeown, K., and Kaiser, G. E. (1988). Automated tutoring in interactive environ-
718 ments: a task centered approach. Technical report, Columbia University.
- 719 [64] Wulff, S. S. and Wulff, D. H. (2004). “of course i’m communicating; I lecture every day”:
720 enhancing teaching and learning in introductory statistics. *Communication Education*, 53(1):92–
721 103.
- 722 [65] Young, J. R. (2012). Inside the Coursera contract: how an upstart company might profit from
723 free courses. *The Chronicle of Higher Education*, 19(7):1–4.
- 724 [66] Ziman, K., Heusser, A. C., Fitzpatrick, P. C., Field, C. E., and Manning, J. R. (2018). Is
725 automatic speech-to-text transcription ready for use in psychological experiments? *Behavior*
726 *Research Methods*, 50:2597–2605.