

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

<sup>4</sup> Paxton C. Fitzpatrick<sup>1</sup>, Andrew C. Heusser<sup>1, 2</sup>, and Jeremy R. Manning<sup>1,\*</sup>

<sup>1</sup>Department of Psychological and Brain Sciences

Dartmouth College, Hanover, NH 03755, USA

<sup>2</sup>Akili Interactive Labs

Boston, MA 02110, USA

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

<sup>5</sup> **Abstract**

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

<sup>17</sup> **Keywords:** education, learning, knowledge, concepts, natural language processing

<sup>18</sup> **Introduction**

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

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

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

45 acquired information to the scaffolding of one's existing knowledge or experience [? ? ? ? ? ]?  
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 [? ]? 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 [? ? ? ].

50 The difference between "understanding" and "memorizing," as framed by researchers in edu-  
51 cation, cognitive psychology, and cognitive neuroscience (e.g., ? ? ? ? ? ), has profound analogs  
52 in the fields of natural language processing and natural language understanding. For example,  
53 considering the raw contents of a document (e.g., its constituent symbols, letters, and words) might  
54 provide some clues as to what the document is about, just as memorizing a passage might provide  
55 some ability to answer simple questions about it. However, text embedding models (e.g., ? ? ? ? ?  
56 ? ? ) also attempt to capture the deeper meaning *underlying* those atomic elements. These models  
57 consider not only the co-occurrences of those elements within and across documents, but also pat-  
58 terns in how those elements appear across different scales (e.g., sentences, paragraphs, chapters,  
59 etc.), the temporal and grammatical properties of the elements, and other high-level characteristics  
60 of how they are used [? ? ]. According to these models, the deep conceptual meaning of a document  
61 may be captured by a feature vector in a high-dimensional representation space, wherein nearby  
62 vectors reflect conceptually related documents. A model that succeeds at capturing an analogue of  
63 "understanding" is able to assign nearby feature vectors to two conceptually related documents,  
64 *even when the specific words contained in those documents have very little overlap.*

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

73 should change accordingly. Learning new concepts should both update our characterizations of  
74 “what is known” and also unlock any now-satisfied dependencies of those newly learned concepts  
75 so that they are “tagged” as available for future learning.

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

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

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

105 **Results**

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

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

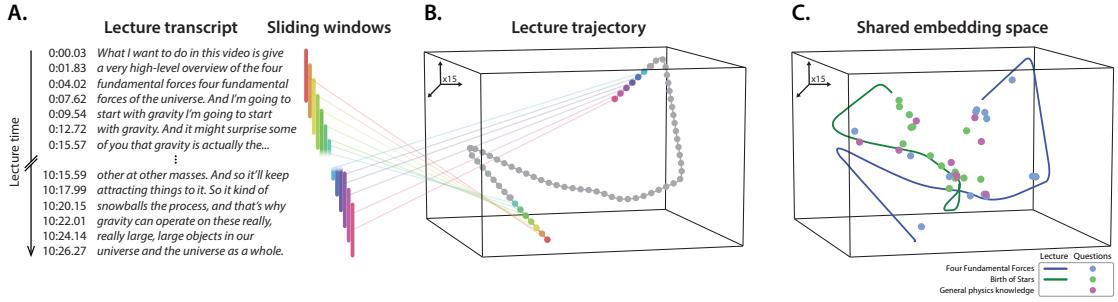


**Figure 1: Experimental paradigm.** Participants alternate between completing three 13-question multiple-choice quizzes and watching two Khan Academy lectures. Each quiz contains a mix of 5 questions about Lecture 1, 5 questions about Lecture 2, and 3 questions about general physics knowledge. The specific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

on 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 the two 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 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 concepts not specifically presented in either video (see Supp. Tab. 1 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 after watching the *Birth of Stars* video (i.e., Lecture 2).

To study in detail how participants’ conceptual knowledge changed over the course of the



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

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

<sup>160</sup> transcript's overlapping text snippets across sliding windows) forms a *trajectory* that captures how  
<sup>161</sup> its conceptual content unfolds over time (Fig. 2B). We resampled these trajectories to a resolution  
<sup>162</sup> of one topic vector for each second of video (i.e., 1 Hz).

<sup>163</sup> We hypothesized that a topic model trained on transcripts of the two lectures should also capture  
<sup>164</sup> the conceptual knowledge probed by each quiz question. If indeed the topic model could capture  
<sup>165</sup> information about the deeper conceptual content of the lectures (i.e., beyond surface-level details  
<sup>166</sup> such as particular word choices), then we should be able to recover a correspondence between each  
<sup>167</sup> lecture and questions *about* each lecture. Importantly, such a correspondence could not solely arise  
<sup>168</sup> from superficial text matching between lecture transcripts and questions, since the lectures and  
<sup>169</sup> questions used different words. Simply comparing the average topic weights from each lecture and  
<sup>170</sup> question set (averaging across time and questions, respectively) reveals a striking correspondence  
<sup>171</sup> (Supp. Fig. 2). Specifically, the average topic weights from Lecture 1 are strongly correlated with the  
<sup>172</sup> average topic weights from Lecture 1 questions ( $r(13) = 0.809, p < 0.001$ , 95% confidence interval  
<sup>173</sup> (CI) = [0.633, 0.962]), and the average topic weights from Lecture 2 are strongly correlated with the  
<sup>174</sup> average topic weights from Lecture 2 questions ( $r(13) = 0.728, p = 0.002$ , 95% CI = [0.456, 0.920]).  
<sup>175</sup> At the same time, the average topic weights from the two lectures are *negatively* correlated with  
<sup>176</sup> their non-matching question sets (Lecture 1 video vs. Lecture 2 questions:  $r(13) = -0.547, p = 0.035$ ,  
<sup>177</sup> 95% CI = [-0.812, -0.231]; Lecture 2 video vs. Lecture 1 questions:  $r(13) = -0.612, p = 0.015$ , 95%  
<sup>178</sup> CI = [-0.874, -0.281]), indicating that the topic model also exhibits some degree of specificity. The  
<sup>179</sup> full set of pairwise comparisons between average topic weights for the lectures and question sets  
<sup>180</sup> is reported in Supplementary Figure 2.

<sup>181</sup> Another, more sensitive, way of summarizing the conceptual content of the lectures and ques-  
<sup>182</sup> tions is to look at *variability* in how topics are weighted over time and across different questions  
<sup>183</sup> (Fig. 3). Intuitively, the variability in the expression of a given topic relates to how much “in-  
<sup>184</sup> formation” [?] the lecture (or question set) reflects about that topic. For example, suppose a  
<sup>185</sup> given topic is weighted on heavily throughout a lecture. That topic might be characteristic of some  
<sup>186</sup> aspect or property of the lecture *overall* (conceptual or otherwise), but unless the topic’s weights  
<sup>187</sup> changed in meaningful ways over time, the topic would be a poor indicator of any *specific* concep-



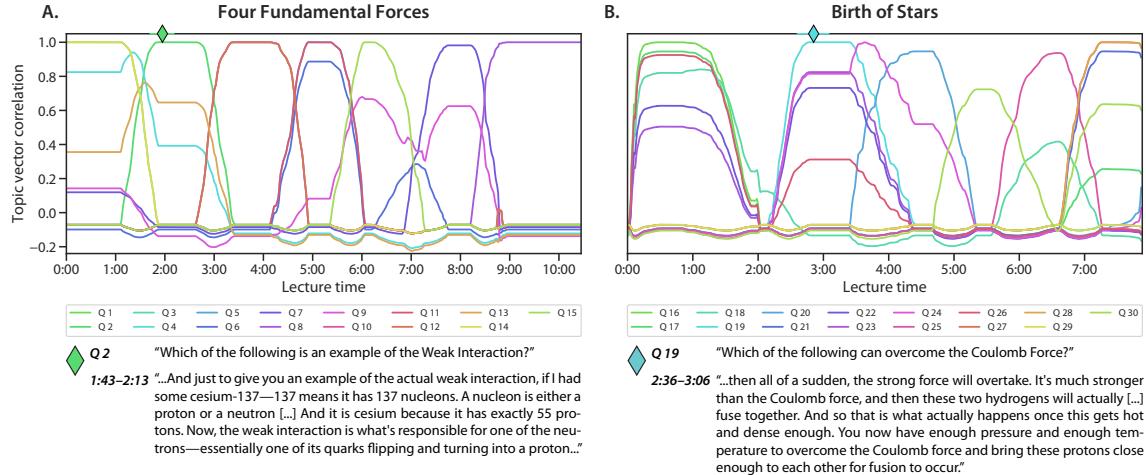
**Figure 3: Lecture and question topic overlap. A. Topic weight variability.** The bar plots display the variance of each topic's weight across lecture timepoints (top row) and questions (bottom row); colors denote topics. The top-weighted words from the most “expressive” (i.e., variable across observations)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

tual content in the lecture. We therefore also compared the variances in topic weights (across time or questions) between the lectures and questions. The variability in topic expression (over time and across questions) was similar for the Lecture 1 video and questions ( $r(13) = 0.824, p < 0.001, 95\% \text{ CI} = [0.696, 0.973]$ ) and the Lecture 2 video and questions ( $r(13) = 0.801, p < 0.001, 95\% \text{ CI} = [0.539, 0.958]$ ). Simultaneously, as reported in Figure 3B, the variability in topic expression across *different* videos and lecture-specific questions (i.e., Lecture 1 video vs. Lecture 2 questions; Lecture 2 video vs. Lecture 1 questions) were negatively correlated, and neither video’s topic variability was reliably correlated with the topic variability across general physics knowledge questions. Taken together, the analyses reported in Figure 3 and Supplementary Figure 2 indicate that a topic model fit to the videos’ transcripts can also reveal correspondences (at a coarse scale) between the lectures and questions.

While an individual lecture may be organized around a single broad theme at a coarse scale, at a finer scale, each moment of a lecture typically covers a narrower range of content. Given the correspondence we found between the variability in topic expression across moments of each

lecture and questions from its corresponding set (Fig. 3), we wondered whether the text embedding model might additionally capture these conceptual relationships at a finer scale. For example, if a particular question asks about the content from one small part of a lecture, we wondered whether the text embeddings could be used to automatically identify the “matching” moment(s) in the lecture. To explore this, we computed the correlation between each question’s topic weights and the topic weights for each second of its corresponding lecture, and found that each question appeared to be temporally specific (Fig. 4). In particular, most questions’ topic vectors were maximally correlated with a well-defined (and relatively narrow) range of timepoints from their corresponding lectures, and the correlations fell off sharply outside of that range. We also qualitatively examined the best-matching intervals for each question by comparing the question’s text 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 probed and the content covered by the best-matching moments of the lectures. Two representative examples are shown at the bottom of Figure 4.

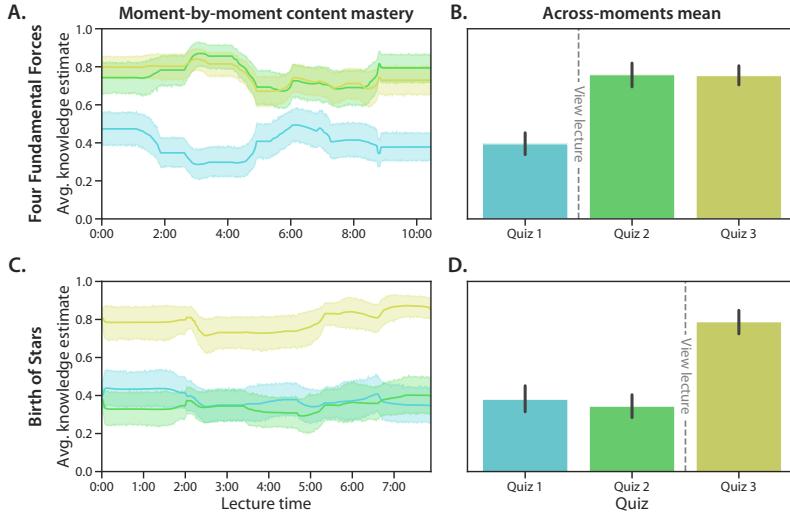
The ability to quantify how much each question is “asking about” the content from each moment of the lectures could enable high-resolution insights into participants’ knowledge. Traditional approaches to estimating how much a student “knows” about the content of a given lecture entail computing the proportion 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 help these two students fill in the “gaps” in their understandings, we might do well to focus specifically on concept *A* for the first student, but to also add in materials pertaining to concepts *B* and *C* for the second student. In other words, raw “proportion-correct” measures may capture *how much* a student knows, but not *what* they know. We wondered whether our modeling framework might enable us to (formally and automatically)



**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 infer participants’ knowledge at the scale of individual concepts (e.g., as captured by a single  
231 moment of a lecture).

232 We developed a simple formula (Eqn. 1) for using a participant’s responses to a small set of  
233 multiple-choice questions to estimate how much the participant “knows” about the concept re-  
234 flected by any arbitrary coordinate,  $x$ , in text embedding space (e.g., the content reflected by any  
235 moment in a lecture they had watched; see *Estimating dynamic knowledge traces*). Essentially, the  
236 estimated knowledge at coordinate  $x$  is given by the weighted average proportion of quiz questions  
237 the participant answered correctly, where the weights reflect how much each question is “about” the  
238 content at  $x$ . When we apply this approach to estimate the participant’s knowledge about the con-  
239 tent presented in each moment of each lecture, we can obtain a detailed timecourse describing how  
240 much “knowledge” the participant has about any part of the lecture. As shown in Figure 5A and C,  
241 we can apply this approach separately for the questions from each quiz participants took through-  
242 out the experiment. From just a few questions per quiz (see *Estimating dynamic knowledge traces*),



**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 a single 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.

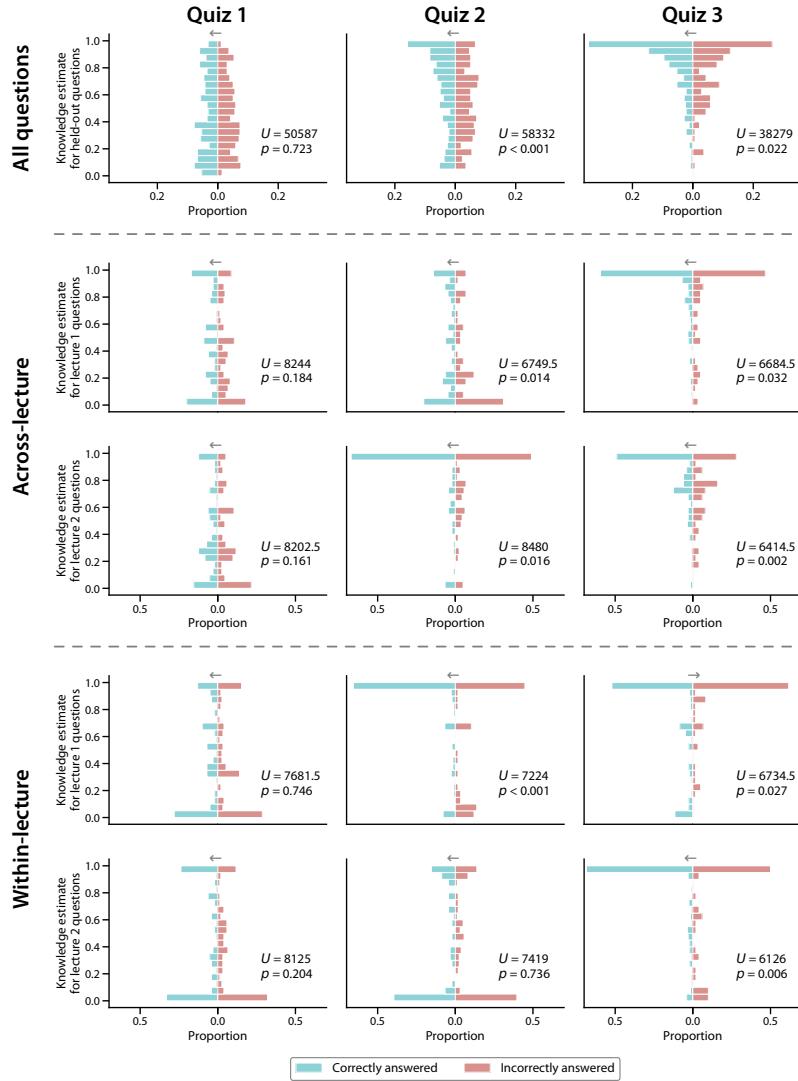
243 we obtain a high-resolution snapshot (at the time each quiz was taken) of what the participants  
 244 knew about any moment’s content, from either of the two lectures they watched (comprising a  
 245 total of 1,100 samples across the two lectures).

246 While the timecourses in Figure 5A and C provide detailed *estimates* about participants’ knowl-  
 247 edge, these estimates are of course only *useful* to the extent that they accurately reflect what  
 248 participants actually know. As one sanity check, we anticipated that the knowledge estimates  
 249 should reflect a content-specific “boost” in participants’ knowledge after watching each lecture.  
 250 In other words, if participants learn about each lecture’s content when they watch each lecture,  
 251 the knowledge estimates should capture that. After watching the *Four Fundamental Forces* lecture,  
 252 participants should exhibit more knowledge for the content of that lecture than they had before,  
 253 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 greater estimated knowledge about the content of the *Birth of Stars* lecture after (versus before) watching it (Fig. 5D). Specifically, since participants watched that lecture after taking Quiz 2 (but before Quiz 3), we hypothesized that their knowledge estimates 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 that question correctly or incorrectly. We developed a statistical approach to test this claim. For each question, in turn, we used Equation 1 to estimate each participant's knowledge at the given question's embedding space coordinate, using all *other* questions that participant answered on the same quiz. For each quiz, we grouped these estimates into two distributions: one for the estimated knowledge at the coordinates of *correctly* answered questions, and another for the estimated knowledge at the coordinates of *incorrectly* answered questions (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

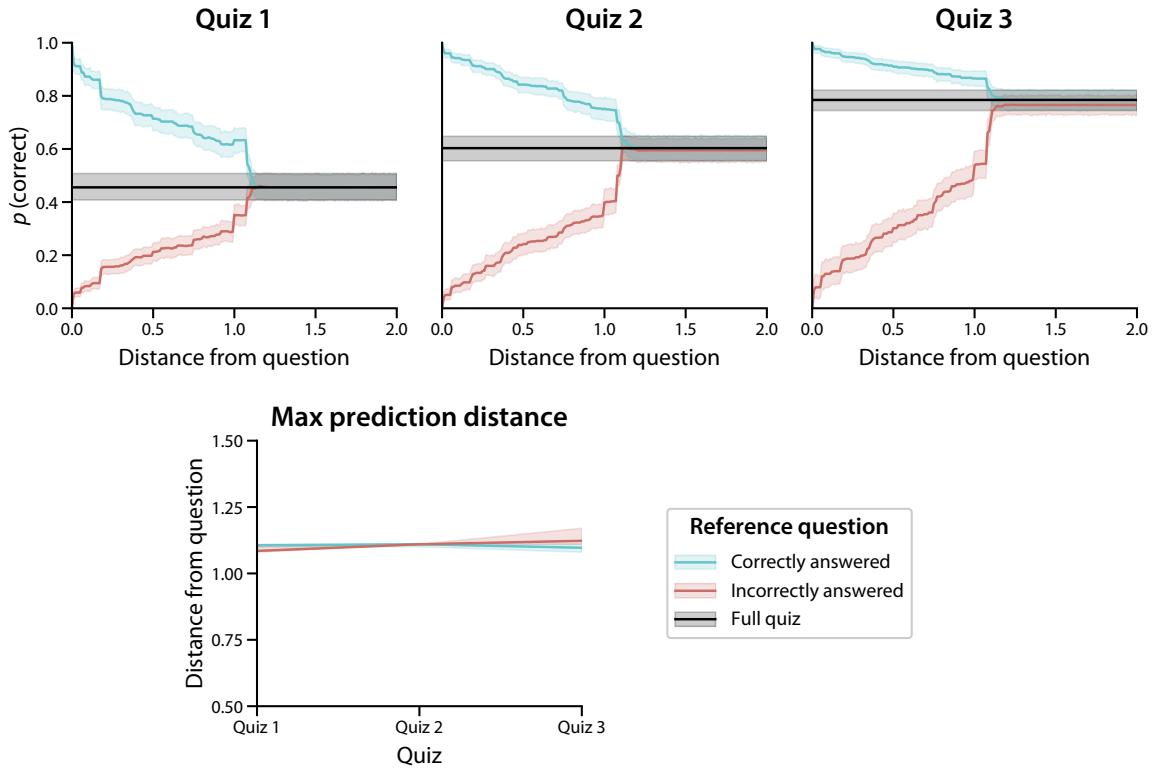


**Figure 6: Estimating knowledge at the embedding coordinates of held-out questions.** Separately for each quiz (panel), we plot the distributions of predicted knowledge at the embedding coordinates of each held-out correctly (blue) or incorrectly (red) answered question. The  $t$ -tests reported in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

282 correctly versus incorrectly answered ( $t(633) = 0.577$ ,  $p = 0.564$ ). This reflects a floor effect: when  
283 knowledge is low everywhere, there is little signal to differentiate between what is known versus  
284 unknown. After watching the first lecture, estimated knowledge for held-out correctly answered  
285 questions (from the second quiz; Fig. 6, middle panel) exhibited a positive shift relative to held-  
286 out incorrectly answered questions ( $t(633) = 3.961$ ,  $p < 0.001$ ). This second quiz provides the  
287 maximally sensitive test for our knowledge predictions, since (if knowledge is estimated accurately)  
288 participants' Quiz 2 responses should demonstrate specific knowledge about Lecture 1 content,  
289 but knowledge about Lecture 2 and general physics concepts should be roughly unchanged from  
290 before they watched Lecture 1. After watching the second lecture, estimated knowledge (from the  
291 third quiz; Fig. 6, right panel) for *all* questions exhibited a positive shift. However, the estimated  
292 knowledge for held-out correctly answered questions remained greater than that for held-out  
293 incorrectly answered questions ( $t(628) = 2.045$ ,  $p = 0.041$ ). This third contrast reflects a ceiling  
294 effect: when knowledge is relatively high everywhere, the signal differentiating what is known  
295 versus unknown is relatively weak. Taken together, this set of analyses demonstrates that our  
296 knowledge prediction framework is most informative when participants exhibit variability in their  
297 knowledge of the content captured by the text embedding model.

298 Knowledge estimates need not be limited to the content of the lectures. As illustrated in  
299 Figure 8, our general approach to estimating knowledge from a small number of quiz questions  
300 may be extended to *any* content, given its text embedding coordinate. To visualize how knowledge  
301 "spreads" through text embedding space to content beyond the lectures participants watched, we  
302 first fit a new topic model to the lectures' sliding windows with (up to)  $k = 100$  topics. Conceptually,  
303 increasing the number of topics used by the model functions to increase the "resolution" of the  
304 embedding space, providing a greater ability to estimate knowledge for content that is highly  
305 similar to (but not precisely the same as) that contained in the two lectures. Aside from increasing  
306 the number of topics from 15 to 100, all other procedures and model parameters were carried over  
307 from the preceding analyses. As in our other analyses, we resampled each lecture's topic trajectory  
308 to 1 Hz and projected each question into a shared text embedding space.

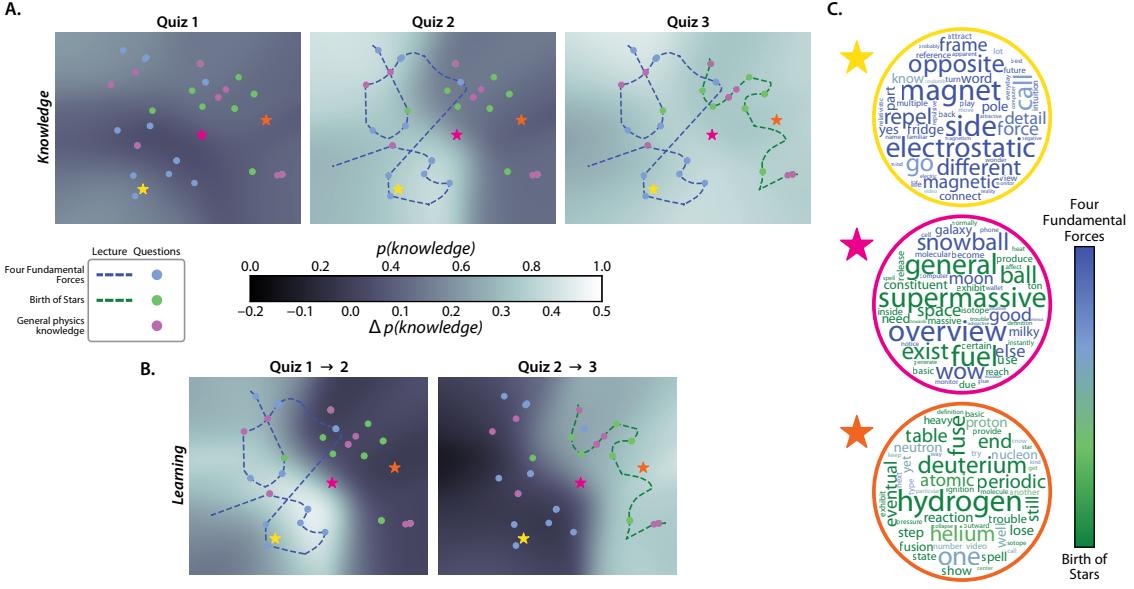
309 We projected the resulting 100-dimensional topic vectors (for each second of video and each quiz



**Figure 7: Caption title.** Caption content.

question) onto a shared 2-dimensional plane (see *Creating knowledge and learning map visualizations*). Next, we sampled points from a  $100 \times 100$  grid of coordinates that evenly tiled a rectangle enclosing the 2D projections of the videos and questions. We used Equation 4 to estimate participants' knowledge at each of these 10,000 sampled locations, and averaged these estimates across participants to obtain an estimated average *knowledge map* (Fig. 8A). Intuitively, the knowledge map constructed from a given quiz's responses provides a visualization of how "much" participants knew about any content expressible by the fitted text embedding model at the point in time when they completed that quiz.

Several features of the resulting knowledge maps are worth noting. The average knowledge map estimated from Quiz 1 responses (Fig. 8A, leftmost map) shows that participants tended to have relatively little knowledge about any parts of the text embedding space (i.e., the shading is



**Figure 8: Mapping out the geometry of knowledge and learning.** **A. Average “knowledge maps” estimated using each quiz.** Each map displays a 2D projection of the estimated knowledge about the content reflected by *all* regions of topic space (see *Creating knowledge and learning map visualizations*). The topic trajectories of the two lectures are indicated by dotted lines (blue: Lecture 1; green: Lecture 2), and the coordinates of each question are indicated by dots (light blue: Lecture 1-related; light green: Lecture 2-related; purple: general physics knowledge). Each map reflects an average across all participants. For individual participants’ maps, see Supplementary Figures 3, 4, and 5. **B. Average “learning maps” estimated between each successive pair of quizzes.** The learning maps follow the same general format as the knowledge maps in Panel A, but here the shading at each coordinate indicates the *difference* between the corresponding coordinates in the indicated pair of knowledge maps—i.e., how much the estimated knowledge “changed” between the two quizzes. Each map reflects an average across all participants. For individual participants’ maps, see Supplementary Figures 6 and 7. **C. Word clouds for sampled points in topic space.** Each word cloud displays the weighted blend of words underlying the topic proportions represented at the corresponding colored star’s location on the maps. In each word cloud, the words’ relative sizes correspond to their relative weights at the starred location, and their colors indicate their relative weights in the *Four Fundamental Forces* (blue) versus *Birth of Stars* (green) lectures, on average, across all timepoints’ topic vectors.

321 relatively dark everywhere). The knowledge map estimated from Quiz 2 responses shows a marked  
322 increase in knowledge on the left side of the map (around roughly the same range of coordinates  
323 traversed by the *Four Fundamental Forces* lecture, indicated by the dotted blue line). In other words,  
324 participants' estimated increase in knowledge is localized to conceptual content that is nearby (i.e.,  
325 related to) the content from the lecture they watched prior to taking Quiz 2. This localization is  
326 non-trivial: these knowledge estimates are informed only by the embedded coordinates of the  
327 *quiz questions*, not by the embeddings of either lecture (see Eqn. 4). Finally, the knowledge map  
328 estimated from Quiz 3 responses shows a second increase in knowledge, localized to the region  
329 surrounding the embedding of the *Birth of Stars* lecture participants watched immediately prior to  
330 taking Quiz 3.

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

337 Because the 2D projection we used to construct the knowledge and learning maps is invertible,  
338 we may gain additional insights into these maps' meaning by reconstructing the original high-  
339 dimensional topic vector for any location on the map we are interested in. For example, this could  
340 serve as a useful tool for an instructor looking to better understand which content areas a student  
341 (or a group of students) knows well (or poorly). As a demonstration, we show the top-weighted  
342 words from the blends of topics reconstructed from three example locations on the maps (Fig. 8C):  
343 one point near the *Four Fundamental Forces* embedding (yellow), a second point near the *Birth of*  
344 *Stars* embedding (orange), and a third point between the two lectures' embeddings (pink). As  
345 shown in the word clouds in the panel, the top-weighted words at the example coordinate near the  
346 *Four Fundamental Forces* embedding tended to be weighted more heavily by the topics expressed  
347 in that lecture. Similarly, the top-weighted words at the example coordinate near the *Birth of Stars*  
348 embedding tended to be weighted more heavily by the topics expressed in *that* lecture. And the

349 top-weighted words at the example coordinate between the two lectures' embeddings show a  
350 roughly even mix of words most strongly associated with each lecture.

## 351 Discussion

352 We developed a computational framework that uses short multiple-choice quizzes to gain nuanced  
353 insights into what learners know and how their knowledge changes with training. First, we show  
354 that our approach can automatically match the conceptual knowledge probed by individual quiz  
355 questions to the corresponding moments in lecture videos when those concepts were presented  
356 (Fig. 4). Next, we demonstrate how we can estimate moment-by-moment "knowledge traces"  
357 that reflect the degree of knowledge participants have about each video's time-varying content,  
358 and capture temporally specific increases in knowledge after viewing each lecture (Fig. 5). We  
359 also show that these knowledge estimates can generalize to held-out questions (Fig. 6). Finally,  
360 we use our framework to construct visual maps that provide snapshot estimates of how much  
361 participants know about any concept within the scope of our text embedding model, and how  
362 much their knowledge of those concepts changes with training (Fig. 8).

363 Over the past several years, the global pandemic has forced many educators to suddenly  
364 adapt to teaching remotely [? ? ? ? ]. This change in world circumstances is happening alongside  
365 (and perhaps accelerating) geometric growth in the availability of high-quality online courses from  
366 platforms such as Khan Academy [? ], Coursera [? ], EdX [? ], and others [? ]. Continued expansion  
367 of the global internet backbone and improvements in computing hardware have also facilitated  
368 improvements in video streaming, enabling videos to be easily shared and viewed by increasingly  
369 large segments of the world's population. This exciting time for online course instruction provides  
370 an opportunity to re-evaluate how we, as a global community, educate ourselves and each other.  
371 For example, we can ask: what defines an effective course or training program? Which aspects of  
372 teaching might be optimized and/or augmented by automated tools? How and why do learning  
373 needs and goals vary across people? How might we lower barriers of access to a high-quality  
374 education?

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

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

Ultimately, our work suggests a rich new line of questions about the geometric “form” of knowledge, how knowledge changes over time, and how we might map out the full space of what an individual knows. Our finding that detailed estimates about knowledge may be obtained from short quizzes shows one way that traditional approaches to evaluation in education may be extended. We hope that these advances might help pave the way for new approaches to teaching

403 or delivering educational content that are tailored to individual students' learning needs and goals.

404 **Materials and methods**

405 **Participants**

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

411 Participants' ages ranged from 18 to 22 years (mean: 19.52 years; standard deviation: 1.09  
412 years). A total of 15 participants reported their gender as male and 35 participants reported their  
413 gender as female. A total of 49 participants reported their native language as "English" and 1  
414 reported having another native language. A total of 47 participants reported their ethnicity as  
415 "Not Hispanic or Latino" and three reported their ethnicity as "Hispanic or Latino." Participants  
416 reported their races as White (32 participants), Asian (14 participants), Black or African American  
417 (5 participants), American Indian or Alaska Native (1 participant), and Native Hawaiian or Other  
418 Pacific Islander (1 participant). (Note that some participants selected multiple racial categories.)

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

426 No participants reported that their focus was currently impaired (e.g., by drugs or alcohol).  
427 Participants reported their current level of alertness, and we converted their responses to numerical

428 scores as follows: “very sluggish” (-2), “a little sluggish” (-1), “neutral” (0), “fairly alert” (1), and  
429 “very alert” (2). Across all participants, a range of alertness levels were reported (range: -2–1;  
430 mean: -0.10; standard deviation: 0.84).

431 Participants reported their undergraduate major(s) as “social sciences” (28 participants), “nat-  
432 ural sciences” (16 participants), “professional” (e.g., pre-med or pre-law; 8 participants), “mathe-  
433 matics and engineering” (7 participants), “humanities” (4 participants), or “undecided” (3 partici-  
434 pants). Note that some participants selected multiple categories for their undergraduate major(s).  
435 We also asked participants about the courses they had taken. In total, 45 participants reported hav-  
436 ing taken at least one Khan Academy course in the past, and 5 reported not having taken any Khan  
437 Academy courses. Of those who reported having watched at least one Khan Academy course,  
438 7 participants reported having watched 1–2 courses, 11 reported having watched 3–5 courses, 8  
439 reported having watched 5–10 courses, and 19 reported having watched 10 or more courses. We  
440 also asked participants about the specific courses they had watched, categorized under different  
441 subject areas. In the “Mathematics” area, participants reported having watched videos on AP  
442 Calculus AB (21 participants), Precalculus (17 participants), Algebra 2 (14 participants), AP Cal-  
443 culus BC (12 participants), Trigonometry (11 participants), Algebra 1 (10 participants), Geometry  
444 (8 participants), Pre-algebra (7 participants), Multivariable Calculus (5 participants), Differential  
445 Equations (5 participants), Statistics and Probability (4 participants), AP Statistics (2 participants),  
446 Linear Algebra (2 participants), Early Math (1 participant), Arithmetic (1 participant), and other  
447 videos not listed in our survey (5 participants). In the “Science and engineering” area, participants  
448 reported having watched videos on Chemistry, AP Chemistry, or Organic Chemistry (21 partic-  
449 ipants); Physics, AP Physics I, or AP Physics II (18 participants); Biology, AP Biology; or High  
450 school Biology (15 participants); Health and Medicine (1 participant); or other videos not listed  
451 in our survey (5 participants). We also asked participants whether they had specifically seen the  
452 videos used in our experiment. Of the 45 participants who reported having having taken at least  
453 one Khan Academy course in the past, 44 participants reported that they had not watched the *Four*  
454 *Fundamental Forces* video, and 1 participant reported that they were not sure whether they had  
455 watched it. All participants reported that they had not watched the *Birth of Stars* video. When

456 we asked participants about non-Khan Academy online courses, they reported having watched  
457 or taken courses on Mathematics (15 participants), Science and engineering (11 participants), Test  
458 preparation (9 participants), Economics and finance (3 participants), Arts and humanities (2 partic-  
459 ipants), Computing (2 participants), and other categories not listed in our survey (17 participants).  
460 Finally, we asked participants about in-person courses they had taken in different subject areas.  
461 They reported taking courses in Mathematics (38 participants), Science and engineering (37 par-  
462 ticipants), Arts and humanities (34 participants), Test preparation (27 participants), Economics  
463 and finance (26 participants), Computing (14 participants), College and careers (7 participants), or  
464 other courses not listed in our survey (6 participants).

465 **Experiment**

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

474 Over the course of the experiment, participants completed three 13-question multiple-choice  
475 quizzes: the first before viewing Lecture 1, the second between Lectures 1 and 2, and the third  
476 after viewing Lecture 2 (see Fig. 1). The questions appearing on each quiz, for each participant,  
477 were randomly chosen from the full set of 39, with the constraints that (a) each quiz contain  
478 exactly 5 questions about Lecture 1, 5 questions about Lecture 2, and 3 questions about general  
479 physics knowledge, and (b) each question appear exactly once for each participant. The orders of  
480 questions on each quiz, and the orders of answer options for each question, were also randomized.  
481 Our experimental protocol was approved by the Committee for the Protection of Human Subjects  
482 at Dartmouth College. We used this experiment to develop and test our computational framework

483 for estimating knowledge and learning.

484 **Analysis**

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

486 We adapted an approach we developed in prior work [? ] to embed each moment of the two  
487 lectures and each question in our pool in a common representational space. Briefly, our approach  
488 uses a topic model (Latent Dirichlet Allocation; ? ), trained on a set of documents, to discover a set  
489 of (up to)  $k$  “topics” or “themes.” Formally, each topic is defined as a distribution of weights over  
490 each word in the model’s vocabulary (i.e., the union of all unique words, across all documents,  
491 excluding “stop words.”). Conceptually, each topic is intended to give larger weights to words that  
492 are semantically related or tend to co-occur in the same documents. After fitting a topic model,  
493 each document in the training set, or any *new* document that contains at least some of the words  
494 in the model’s vocabulary, may be represented as a  $k$ -dimensional vector describing how much  
495 the document (most probably) reflects each topic. To select an appropriate  $k$  for our model, we  
496 identified the minimum number of topics that yielded at least one “unused” topic (i.e., in which  
497 all words in the vocabulary were assigned uniform weights) after training. This indicated that  
498 the number of topics was sufficient to capture the set of latent themes present in the two lectures  
499 (from which we constructed our document corpus, as described below). We found this value to  
500 be  $k = 15$  topics. The distribution of weights over words in the vocabulary for each discovered  
501 topic is shown in Supplementary Figure 1, and each topic’s top-weighted words may be found in  
502 Supplementary Table 2.

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

509 of sufficiently near-human quality to yield reliable data in behavioral studies; ? ) when developing  
510 our framework in order to make it more directly extensible and adaptable by others in the future.  
511 We fetched these automated transcripts using the `youtube-transcript-api` Python package [? ]  
512 ]. The transcripts consisted of one timestamped line of text for every few seconds (mean: 2.34 s;  
513 standard deviation: 0.83 s) of spoken content in the video (i.e., corresponding to each individual  
514 caption that would appear on-screen if viewing the lecture via YouTube, and when those lines  
515 would appear). We defined a sliding window length of (up to)  $w = 30$  transcript lines, and  
516 assigned each window a timestamp corresponding to the midpoint between the timestamps for its  
517 first and last lines. These sliding windows ramped up and down in length at the beginning and  
518 end of each transcript, respectively. In other words, each transcript's first sliding window covered  
519 only its first line, the second sliding window covered the first two lines, and so on. This ensured  
520 that each line from the transcripts appeared in the same number ( $w$ ) of sliding windows. After  
521 performing various standard text preprocessing (e.g., normalizing case, lemmatizing, removing  
522 punctuation and stop-words), we treated the text from each sliding window as a single "document,"  
523 and combined these documents across the two videos' windows to create a single training corpus  
524 for the topic model.

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

532 We transformed each sliding window's text into a topic vector, and then used linear interpolation  
533 (independently for each topic dimension) to resample the resulting timeseries to one vector  
534 per second. We also used the fitted model to obtain topic vectors for each question in our pool (see  
535 Supp. Tab. 1). Taken together, we obtained a *trajectory* for each video, describing its path through  
536 topic space, and a single coordinate for each question (Fig. 2C). Embedding both videos and all of

537 the questions using a common model enables us to compare the content from different moments  
 538 of videos, compare the content across videos, and estimate potential associations between specific  
 539 questions and specific moments of video.

540 **Estimating dynamic knowledge traces**

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

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

543 where

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

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

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

558 **Creating knowledge and learning map visualizations**

559 An important feature of our approach is that, given a trained text embedding model and partic-  
560 ipants’ quiz performance on each question, we can estimate their knowledge about *any* content  
561 expressible by the embedding model—not solely the content explicitly probed by the quiz ques-  
562 tions, or even appearing in the lectures. To visualize these estimates (Fig. 8, Supp. Figs. 3, 4, 5, 6,  
563 and 7), we used Uniform Manifold Approximation and Projection (UMAP; ??) to construct a  
564 2D projection of the text embedding space. Sampling the original 100-dimensional space at high  
565 resolution to obtain an adequate set of topic vectors spanning the embedding space would be  
566 computationally intractable. However, sampling a 2D grid is trivial.

567 At a high level, the UMAP algorithm obtains low-dimensional embeddings by minimizing  
568 the cross-entropy between the pairwise (clustered) distances between the observations in their  
569 original (e.g., 100-dimensional) space and the pairwise (clustered) distances in the low-dimensional  
570 embedding space (in our approach, the embedding space is 2D). In our implementation, pairwise  
571 distances in the original high-dimensional space were defined as 1 minus the correlation between  
572 each pair of coordinates, and pairwise distances in the low-dimensional embedding space were  
573 defined as the Euclidean distance between each pair of coordinates.

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

584 After embedding both lectures’ topic trajectories and the topic vectors of every question, we

585 defined a rectangle enclosing the 2D projections of the lectures' and quizzes' embeddings. We then  
586 sampled points from a regular  $100 \times 100$  grid of coordinates that evenly tiled this enclosing rectangle.  
587 We sought to estimate participants' knowledge (and learning, i.e., changes in knowledge) at each  
588 of the resulting 10,000 coordinates.

589 To generate our estimates, we placed a set of 39 radial basis functions (RBFs) throughout the  
590 embedding space, centered on the 2D projections for each question (i.e., we included one RBF for  
591 each question). At coordinate  $x$ , the value of an RBF centered on a question's coordinate  $\mu$ , is given  
592 by:

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

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

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

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

## 600 Author contributions

601 Conceptualization: PCF, ACH, and JRM. Methodology: PCF, ACH, and JRM. Software: PCF.  
602 Validation: PCF. Formal analysis: PCF. Resources: PCF, ACH, and JRM. Data curation: PCF.  
603 Writing (original draft): JRM. Writing (review and editing): PCF, ACH, and JRM. Visualization:  
604 PCF and JRM. Supervision: JRM. Project administration: PCF. Funding acquisition: JRM.

605 **Data and code availability**

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

609 **Acknowledgements**

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