

1            **Geometric models reveal the hidden structure of**  
2            **conceptual knowledge**

3            Paxton C. Fitzpatrick and Jeremy R. Manning\*

Dartmouth College

\*Corresponding author: jeremy.r.manning@dartmouth.edu

4            **Abstract**

5            We develop a mathematical framework, based on natural language processing models, for  
6            tracking and characterizing the acquisition of conceptual knowledge. Our approach embeds each  
7            concept in a high dimensional representation space, where nearby coordinates reflect similar or  
8            related concepts. We tested our approach using behavioral data collected from a group of  
9            college students. In the experiment, we asked the participants to answer sets of quiz questions  
10          interleaved between watching two course videos from the Khan Academy platform. We applied  
11          our framework to the videos' transcripts, and to text of the quiz questions, to quantify the  
12          content of each moment of video and each quiz question. We used these embeddings, along with  
13          participants' quiz responses, to track how the learners' knowledge changed after watching each  
14          video. Our findings show how a limited set of quiz questions may be used to construct rich and  
15          meaningful representations of what each learner knows, and how their knowledge changes over  
16          time as they learn.

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

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

<sup>19</sup> How do we acquire conceptual knowledge? Memorizing course lectures or textbook chapters by  
<sup>20</sup> rote can lead to the superficial *appearance* of understanding the underlying content, but achieving  
<sup>21</sup> true conceptual understanding seems to require something deeper and richer. Does conceptual  
<sup>22</sup> understanding entail connecting newly acquired information to the scaffolding of one's existing  
<sup>23</sup> knowledge or experience [1, 5, 7, 8, 22]? Or weaving a lecture's atomic elements (e.g., its compo-  
<sup>24</sup> nent words) into a structured network that describes how those individual elements are related?  
<sup>25</sup> Conceptual understanding could also involve building a mental model that transcends the mean-  
<sup>26</sup> ings of those individual atomic elements by reflecting the deeper meaning underlying the gestalt  
<sup>27</sup> whole [14, 16, 21].

<sup>28</sup> The difference between “understanding” and “memorizing,” as framed by the researchers  
<sup>29</sup> in education, cognitive psychology, and cognitive neuroscience [9, 10, 13, 16, 21] has profound  
<sup>30</sup> analogs in the fields of natural language processing and natural language understanding. For  
<sup>31</sup> example, considering the raw contents of a document (e.g., its constituent symbols, letters, and  
<sup>32</sup> words) might provide some information about what the document is about, just as memorizing  
<sup>33</sup> a passage might be used to answer simple questions about the passage [e.g., whether it might  
<sup>34</sup> contain words related to furniture versus physics; 2, 3, 15]. However, modern natural language  
<sup>35</sup> processing models [e.g., 4, 6, 20] also attempt to capture the deeper meaning *underlying* those  
<sup>36</sup> atomic elements. These models consider not only the co-occurrences of those elements within  
<sup>37</sup> and across documents, but also patterns in how those elements appear across different scales (e.g.,  
<sup>38</sup> sentences, paragraphs, chapters, etc.), the temporal and grammatical properties of the elements,  
<sup>39</sup> and other high-level characteristics of how they are used [17, 18]. According to these models, the  
<sup>40</sup> deep conceptual meaning of a document may be captured by a feature vector in a high-dimensional  
<sup>41</sup> representation space, where nearby vectors reflect conceptually related documents. A model that  
<sup>42</sup> succeeds at capturing an analog of “understanding” is able to assign nearby feature vectors to  
<sup>43</sup> two conceptually related documents, *even when the words contained in those documents have very little  
44 overlap.*

45        What form might the representation of the sum total of a person’s knowledge take? First,  
46    we might require a means of systematically describing or representing the nearly infinite set of  
47    possible things a person could know. Second, we might want to account for potential associations  
48    between different concepts. For example, the concepts of “fish” and “water” might be associated in  
49    the sense that fish live in water. Third, knowledge may have a critical dependency structure, such  
50    that knowing about a particular concept might require first knowing about a set of other concepts.  
51    For example, understanding the concept of a fish swimming in water first requires understanding  
52    what fish and water *are*. Fourth, as we learn, our “current state of knowledge” should change  
53    accordingly. Learning new concepts should both update our characterizations of “what is known”  
54    and should also unlock any now-satisfied dependencies of that newly learned concept so that they  
55    are “tagged” as available for future learning.

56        Here we develop a framework for modelling how knowledge is acquired during learning. The  
57    central idea is to use text embedding models to define the coordinate systems of two maps: (a) a  
58    *knowledge map* that describes the extent to which each concept is currently known and (b) a *learning*  
59    *map* that describes the extent to which each concept could be learned. Each location on these maps  
60    represents a single concept, and the geometries are defined such that related concepts are located  
61    nearby in space. We use this framework to analyzing and interpreting behavioral data collected  
62    from an experiment that has participants watch and answer conceptual questions about a series of  
63    recorded course lectures.

64        Our primary research goal is to advance our understanding of what it means to acquire deep  
65    real-world conceptual knowledge. Traditional laboratory approaches to studying learning and  
66    memory (e.g., list learning studies) often draw little distinction between memorization and under-  
67    standing. Instead, these studies typically focus on whether information is effectively encoded or  
68    retrieved, rather than whether the information is *understood*. Approaches to studying conceptual  
69    learning, such as category learning experiments, can start to investigate the distinction between  
70    memorization and understanding, often by training participants to distinguish arbitrary or ran-  
71    dom features in otherwise meaningless categorized stimuli. However the objective of real-world  
72    training, or learning from life experiences more generally, is often to develop new knowledge

73 that may be applied in *useful* ways in the future. In this sense, the gap between modern learning  
74 theories and modern pedagogical approaches and classroom learning strategies is enormous: most  
75 of our theories about *how* people learn are inspired by experimental paradigms and models that  
76 have only peripheral relevance to the kinds of learning that students and teachers actually seek.  
77 To help bridge this gap, our study uses course materials from real online courses to inform, fit, and  
78 test models of real-world conceptual learning.

## 79 **Results**

## 80 **Discussion**

## 81 **Materials and methods**

### 82 **Participants**

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

88 Participants' ages ranged from 18 to 22 years (mean: 19.52 years; standard deviation: 1.09  
89 years). A total of 15 participants reported their gender as male and 35 participants reported their  
90 gender as female. A total of 49 participants reported their native language as "English" and 1  
91 reported having another native language. A total of 47 participants reported their ethnicity as  
92 "Not Hispanic or Latino" and three reported their ethnicity as "Hispanic or Latino." Participants  
93 reported their races as White (32 participants), Asian (14 participants), Black or African American  
94 (5 participants), American Indian or Alaska Native (1 participant), and Native Hawaiian or Other  
95 Pacific Islander (1 participant). (Note that some participants selected multiple racial categories.)

96 A total of 49 participants reporting having normal hearing and 1 participant reported having  
97 some hearing impairment. A total of 49 participants reported having normal color vision and 1  
98 participant reported being color blind. Participants reported having had, on the night prior to  
99 testing, 2 – 4 hours of sleep (1 participant), 4 – 6 hours of sleep (9 participants), 6 – 8 hours of sleep  
100 (35 participants), or 8+ hours of sleep (5 participants). They reported having consumed, on the  
101 same day and leading up to their testing session, 0 cups of coffee (38 participants), 1 cup of coffee  
102 (10 participants), 3 cups of coffee (1 participant), or 4+ cups of coffee (1 participant).

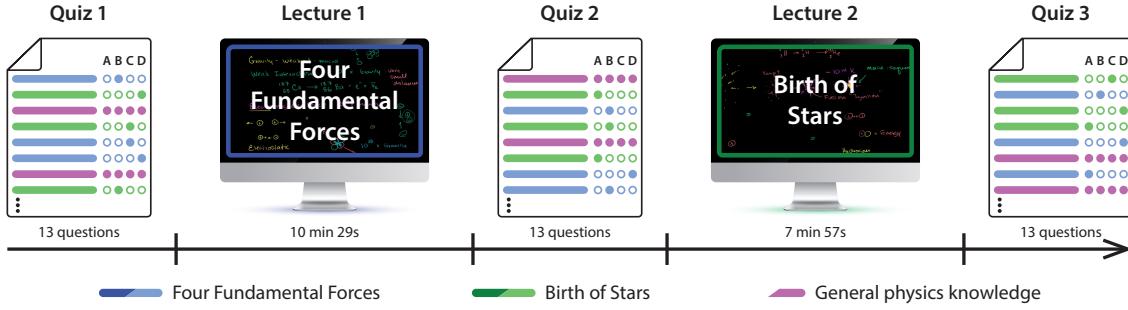
103 No participants reported that their focus was currently impaired (e.g., by drugs or alcohol).  
104 Participants reported their current level of alertness, and we converted their responses to numerical  
105 scores as follows: “very sluggish” (-2), “a little sluggish” (-1), “neutral” (0), “fairly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

108 Participants reported their undergraduate major(s) as Social Sciences (28 participants), Natural  
109 sciences (16), Professional (e.g., pre-med or pre-law; 8 participants), Mathematics and engineering  
110 (7 participants), Humanities (4 participants), or Undecided (3 participants). Note that some par-  
111 ticipants selected multiple categories for their undergraduate major. We also asked participants  
112 about the courses they had taken. In total, 46 participants reported having taken at least one Khan  
113 academy course in the past or being familiar with the Khan academy, and 4 reported not having  
114 taken any Khan academy courses. Of the participants who reported having watched at least one  
115 Khan academy course, 1 participant declined to report the number of courses they had watched;  
116 7 participants reported having watched 1–2 courses; 11 reported having watched 3–5 courses; 8  
117 reported having watched 5–10 courses; and 19 reported having watched 10 or more courses. We  
118 also asked participants about the specific courses they had watched, categorized under different  
119 subject areas. In the “Mathematics” area participants reported having watched videos on AP  
120 Calculus AB (21 participants), Precalculus (17 participants), Algebra 2 (14 participants), AP Cal-  
121 culus BC (12 participants), Trigonometry (11 participants), Algebra 1 (10 participants), Geometry  
122 (8 participants), Pre-algebra (7 participants), Multivariable Calculus (5 participants), Differential  
123 Equations (5 participants), Statistics and Probability (4 participants), AP Statistics (2 participants),

<sup>124</sup> Linear Algebra (2 participants), Early Math (1 participant), Arithmetic (1 participant), and other  
<sup>125</sup> videos not listed in our survey (6 participants). In the “Science and engineering” area participants  
<sup>126</sup> reported having watched videos on Chemistry, AP Chemistry, or Organic Chemistry (21 partic-  
<sup>127</sup> ipants); Physics, AP Physics I, or AP Physics II (15 participants); Biology, AP Biology; or High  
<sup>128</sup> school Biology (15 participants); Health and Medicine (1 participant); or other videos not listed in  
<sup>129</sup> our survey (20 participants). We also asked participants if they had specifically seen the videos  
<sup>130</sup> used in our experiment. When we asked about the *Four Fundamental Forces* video, 45 participants  
<sup>131</sup> reported not having watched it before, 1 participant reported that they were not sure if they had  
<sup>132</sup> watched it before, and 4 participants declined to respond. When we asked about the *Birth of*  
<sup>133</sup> *Stars* video, 46 participants reported not having watched it before and 4 participants declined to  
<sup>134</sup> respond. When we asked participants about non-Khan academy online courses, they reported  
<sup>135</sup> having watched or taken courses on Mathematics (15 participants), Science and engineering (11  
<sup>136</sup> participants), Test preparation (9 participants), Economics and finance (3 participants), Arts and  
<sup>137</sup> humanities (2 participants), Computing (2 participants), and other categories not listed in our  
<sup>138</sup> survey (18 participants). Finally, we asked participants about in-person courses they had taken in  
<sup>139</sup> different subject areas. They reported taking courses in Mathematics (39 participants), Science and  
<sup>140</sup> engineering (38 participants), Arts and humanities (35 participants), Test preparation (27 partici-  
<sup>141</sup> pants), Economics and finance (26 participants), Computing (15 participants), College and careers  
<sup>142</sup> (7 participants), or other courses not listed in our survey (6 participants).

## <sup>143</sup> Experiment

<sup>144</sup> We hand-selected two roughly 10-minute course videos from the Khan Academy platform: *The*  
<sup>145</sup> *Four Fundamental Forces* (an introduction to gravity, electromagnetism, the weak nuclear force, and  
<sup>146</sup> the strong nuclear force; duration: 10 minutes and 29 seconds) and *Birth of Stars* (an introduction  
<sup>147</sup> to how stars are formed; duration: 7 minutes and 57 seconds). We hand-wrote 39 multiple  
<sup>148</sup> choice questions: 15 about the conceptual content of *The Four Fundamental Forces*, another 15 about  
<sup>149</sup> the conceptual content of *Birth of Stars*, and 9 other questions that tested for general conceptual  
<sup>150</sup> knowledge about basic physics (covering material that was not presented in either video). The full



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

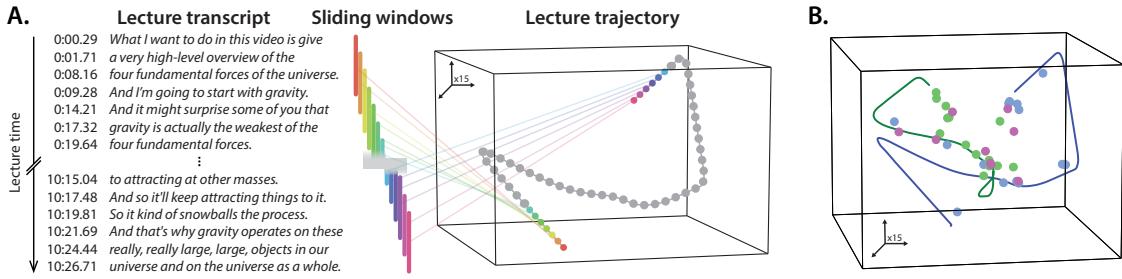
151 set of questions may be found in Table S1.

152 Participants began the main experiment by answering a battery of 13 randomly selected ques-  
153 tions (chosen from the full set of 39). Then they watched the *The Four Fundamental Forces* video.  
154 Next, they answered a second set of 13 questions (chosen at random from the remaining 26 ques-  
155 tions). Fourth, participants watch the *Birth of Stars* video, and finally they answered the remaining  
156 13 questions. Our experimental procedure is diagramed in Figure 1. We used the experiment to  
157 develop and test our computational framework for estimating knowledge and learning maps.

## 158 Analysis

### 159 Constructing text embeddings of multiple videos and questions

160 We extended an approach developed by [12] to construct text embeddings for each moment of each  
161 lecture, and of each question in our pool. Briefly, our approach uses a topic model [3], trained on a  
162 set of documents, to discover a set of  $k$  “topics” or “themes.” Formally, each topic is defined as a set  
163 of weights over each word in the model’s vocabulary (i.e., the union of all unique words, across all  
164 documents, excluding “stop words.”). Conceptually, each topic is intended to give larger weights  
165 to set of words that appear conceptually related or that tend to co-occur in the same documents.  
166 After fitting a topic model, each document in the training set, or any *new* document that contains at  
167 least some of the words in the model’s vocabulary, may be represented as a  $k$ -dimensional vector



**Figure 2: Constructing video content *trajectories*.** **A. Building a document pool from sliding windows**. We decompose each video’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. After training a text embedding model using the two videos’ sliding windows, along with the text from each question in our pool (Tab. S1), we construct “trajectories” through text embedding space by joining the embedding coordinates of successive sliding windows from each video. **B. Embedding multiple videos and questions.** Applying the same text embedding approach to each video, along with the text of each question, results in one trajectory per video and one embedding coordinate (dot) per question (blue: Four Fundamental Forces; green: Birth of Stars; pink: general physics knowledge). Here we have projected the 15-dimensional embeddings into a 3D space using Uniform Manifold Approximation and Projection [UMAP; 19].

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

As illustrated in Figure 2A, we start by building up a corpus of documents using overlapping sliding windows that span each video’s transcript. Khan Academy videos are hosted on the YouTube platform, and all YouTube videos are run through Google’s speech-to-text API [11] to derive a timestamped transcript of any detected speech in the video. The resulting transcripts contain one timestamped row per line, and each line generally corresponds to a few seconds of spoken content from the video. We defined a sliding window length of (up to)  $w = 30$  transcript lines, and we assigned each window a timestamp according to the midpoint between its first and last lines’ timestamps. These sliding windows ramped up and down in length at the very beginning and end of the transcript, respectively. In other words, the first sliding window covered only the first line from the transcript; the second sliding window covered the first two lines; and so on. This ensured that each line of the transcript appeared in the same number ( $w$ ) of sliding windows. We treated the text from each sliding window as a single “document,” and we combined these documents across the two videos’ windows to create a single training corpus for the topic

183 model. The top words from each of the 15 discovered topics may be found in Table S2.

184 After fitting a topic model to each videos' transcripts, we could use the trained model to  
185 transform arbitrary (potentially new) documents into  $k$ -dimensional topic vectors. A convenient  
186 property of these topic vectors is that documents that reflect similar blends of topics (i.e., documents  
187 that reflect similar themes, according to the model) will yield similar (in terms of Euclidean distance,  
188 correlation, etc.) topic vectors. In general, the similarity between different documents' topic vectors  
189 may be used to characterize the similarity in content between the documents.

190 We transformed each sliding window's text into a topic vector, and then used linear interpo-  
191 lation (independently for each topic dimension) to resample the resulting timeseries to once per  
192 second. This yielded a single topic vector for each second of each video. We also used the fitted  
193 model to obtain topic vectors for each question in our pool (Tab. S1). Taken together, we obtained  
194 a *trajectory* for each video, describing its path through topic space, and a single coordinate for each  
195 question (Fig. 2B). Embedding both videos and all of the questions using a common model enables  
196 us to compare the content from different moments of videos, compare the content across videos,  
197 and estimate potential associations between specific questions and specific moments of video.

198 **Estimating dynamic knowledge traces**

199 **Estimating held-out conceptual knowledge**

200 **Creating knowledge and learning map visualizations**

201 **References**

202 [1] Blaye, A., Bernard-Peyron, V., Paour, J.-L., and Bonthoux, F. (2006). Category flexibility in chil-  
203 dren: distinguishing response flexibility from conceptual flexibility; the protracted development  
204 of taxonomic representations. *European Journal of Developmental Psychology*, 3(2):163–188.

205 [2] Blei, D. M. and Lafferty, J. D. (2006). Dynamic topic models. In *Proceedings of the International*  
206 *Conference on Machine Learning*, pages 113–120, New York, NY. Association for Computing  
207 Machinery.

- 208 [3] Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine*  
209 *Learning Research*, 3:993–1022.
- 210 [4] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A.,  
211 Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child,  
212 R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M.,  
213 Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei,  
214 D. (2020). Language models are few-shot learners. *arXiv*, 2005.14165.
- 215 [5] Caramazza, A. and Mahon, B. Z. (2003). The organization of conceptual knowledge: the  
216 evidence from category-specific semantic deficits. *Trends in Cognitive Sciences*, 7(8):354–361.
- 217 [6] Cer, D., Yang, Y., Kong, S. Y., Hua, N., Limtiaco, N., John, R. S., Constant, N., Guajardo-  
218 Cespedes, M., Yuan, S., Tar, C., Sung, Y.-H., Strope, B., and Kurzweil, R. (2018). Universal  
219 sentence encoder. *arXiv*, 1803.11175.
- 220 [7] Constantinescu, A. O., O'Reilly, J. X., and Behrens, T. E. J. (2016). Organizing conceptual  
221 knowledge in humans with a gridlike code. *Science*, 352(6292):1464–1468.
- 222 [8] Deacon, D., Grose-Fifer, J., Yang, C. M., Stanick, V., Hewitt, S., and Dynowska, A. (2004).  
223 Evidence for a new conceptualization of semantic representation in the left and right cerebral  
224 hemispheres. *Cortex*, 40(3):467–478.
- 225 [9] Gallagher, J. J. (2000). Teaching for understanding and application of science knowledge. *School*  
226 *Science and Mathematics*, 100(6):310–318.
- 227 [10] Hall, R. and Greeno, J. (2008). *21st century education: A reference handbook*, chapter Conceptual  
228 learning, pages 212–221. Sage Publications.
- 229 [11] Halpern, Y., Hall, K. B., Schogol, V., Riley, M., Roark, B., Skobeltsyn, G., and Bäuml, M. (2016).  
230 Contextual prediction models for speech recognition. In *Interspeech*, pages 2338–2342.

- 231 [12] Heusser, A. C., Fitzpatrick, P. C., and Manning, J. R. (2021). Geometric models reveal be-  
232 havioral and neural signatures of transforming naturalistic experiences into episodic memories.  
233 *Nature Human Behavior*, 5:905–919.
- 234 [13] Katona, G. (1940). *Organizing and memorizing: studies in the psychology of learning and teaching*.  
235 Columbia University Press.
- 236 [14] Kintsch (1970). *Learning, memory, and conceptual processes*. Wiley.
- 237 [15] Landauer, T. K. and Dumais, S. T. (1997). A solution to Plato’s problem: the latent semantic  
238 analysis theory of acquisition, induction, and representation of knowledge. *Psychological Review*,  
239 104:211–240.
- 240 [16] Maclellan, E. (2005). Conceptual learning: the priority for higher education. *British Journal of  
241 Educational Studies*, 53(2):129–147.
- 242 [17] Manning, J. R. (2020). Context reinstatement. In Kahana, M. J. and Wagner, A. D., editors,  
243 *Handbook of Human Memory*. Oxford University Press.
- 244 [18] Manning, J. R. (2021). Episodic memory: mental time travel or a quantum “memory wave”  
245 function? *Psychological Review*, 128(4):711–725.
- 246 [19] McInnes, L., Healy, J., and Melville, J. (2018). UMAP: Uniform manifold approximation and  
247 projection for dimension reduction. *arXiv*, 1802(03426).
- 248 [20] Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word represen-  
249 tations in vector space. *arXiv*, 1301.3781.
- 250 [21] Scott, P., Asoko, H., and Leach, J. (2007). *Handbook of research on science education*, chapter  
251 Student conceptions and conceptual learning in science. Routledge.
- 252 [22] Simon, M. A., Tzur, R., Heinz, K., and Kinzel, M. (2004). Explicating a mechanism for  
253 conceptual learning: elaborating the construct of reflective abstraction. *Journal for Research in  
254 Mathematics Education*, 35(5):305–329.