

# A Toolbox for Modelling Engagement with Educational Videos

Yuxiang Qiu\*, Karim Djemili\*, Denis Elezi\*, Aaneel Shalman Srazali\*, María Pérez-Ortiz, Emine Yilmaz, John Shawe-Taylor, Sahan Bulathwela

Department of Computer Science, University College London  
Gower Street, London WC1E 6BT, UK  
m.bulathwela@ucl.ac.uk

## Abstract

With the advancement and utility of Artificial Intelligence (AI), personalising education to a global population could be a cornerstone of new educational systems in the future. This work presents the *PEEK*C dataset and the *TrueLearn* Python library, which contains a dataset and a series of online learner state models that are essential to facilitate research on learner engagement modelling. *TrueLearn* family of models was designed following the "open learner" concept, using humanly-intuitive user representations. This family of scalable, online models also help end-users visualise the learner models, which may in the future facilitate user interaction with their models/recommenders. The extensive documentation and coding examples make the library highly accessible to both machine learning developers and educational data mining and learning analytics practitioners. The experiments show the utility of both the dataset and the library with predictive performance significantly exceeding comparative baseline models. The dataset contains a large amount of AI-related educational videos, which are of interest for building and validating AI-specific educational recommenders.

## Introduction

It has been shown that personalised one-on-one learning could lead to improving learning gains by two standard deviations (Bloom 1984). With this goal in sight, and the ambition to democratise education to a world population, we require responsible intelligent systems that can bring scalable, personalised and governable models to a mass of learners (Bulathwela et al. 2024). Intelligent Tutoring Systems (ITS), the go-to solution, is practical for courses with a limited number of learning materials and heavily relies on testing users for knowledge. However, today's world has access to 100,000s of rich educational videos, PDFs and podcasts that can be matched to a global population of lifelong learners. Educational recommenders have the opportunity to leverage implicit interaction signals (such as clicks and watch time) to personalise and support learning for informal lifelong learners (Bulathwela, Sahan and Pérez-Ortiz, María and Yilmaz, Emine and Shawe-Taylor, John 2022). Furthermore, the scarcity of publicly available datasets of learners

in the wild engaging with educational materials is a major deterrent to creating scalable educational recommenders.

The contributions of this paper are two-fold. Firstly, we create and release to the public the **Personalised Educational Engagement linked to Knowledge Components (PEEK**C) dataset, with more than 20,000 informal learners watching educational videos in an in-the-wild setting (i.e. learning informally). The videos in the PEEKC dataset majorly contain concepts related to Artificial Intelligence (AI) and Machine Learning (ML) making it a valuable resource for AI-assisted education on these topics. Secondly, we develop and release *TrueLearn*, an open-source Python library packaging state-of-the-art Bayesian models and visualisation tools for leveraging scalable, online learning, transparent learner models. The library contains different components that will enable i) creating content representations of learning resources ii) managing user/learner states, iii) modelling the state evolution of learners using interactions and iv) evaluating engagement predictions. This work uses the PEEKC dataset to empirically demonstrate the predictive capabilities of the different models within the *TrueLearn* library.

## Related Work

The scarcity of publicly available datasets for predicting learner engagement with educational videos constrains the growth of the personalisation of AI education. PEEKC is the first and largest learner video engagement dataset publicly released with humanly interpretable Wikipedia concepts and the concept coverage associated with the video lecture fragments. Next, we present relevant works to i) the PEEKC dataset (not only related datasets but also research on the approaches that were used to generate it, e.g. Wikification) and ii) our novel machine learning library.

## Related Datasets

*Knowledge Tracing* (Corbett and Anderson 1994; Piech et al. 2015) is the most active research area in the learner modelling domain, focusing on modelling knowledge/skill mastery of learners based on test-taking. ASSISTments data (Selent, Patikorn, and Heffernan 2016), which records learners solving mathematics problems, is often used in literature while this data is mathematics education focused. Additionally, problem-solving interactions (Choi et al. 2020), multiple choice question answering (Wang et al. 2020, 2021)

\*These authors contributed equally.

Copyright © 2024, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

datasets exist publicly while none of them includes implicit feedback related to consuming educational videos. MOOC-Cube dataset contains a spectrum of different statistics relating to learner-MOOC interactions including implicit and explicit test-taking activity (Yu et al. 2020a). Although this dataset may contain data that can be used to predict learner engagement which has been used for course recommendation (Deng et al. 2023), the pre-organisation of courses takes away the in-the-wild choices learners make in choosing videos/fragments to watch while the courses in MOOC-Cube are not limited to AI. On the contrary, PEEKC dataset presents over 20,000 informal learners watching AI-related videos with fragment-level annotation of videos providing more granularity at a time segment/fragment level information retrieval is gaining interest (Yu et al. 2020b).

### Extracting Knowledge Components

ITSs often rely on expert labelling of the **Knowledge Components (KCs)** (Selent, Patikorn, and Heffernan 2016), which is time-consuming and not scalable. Unsupervised learning approaches are also potential candidates. Latent Dirichlet Allocation (LDA) has widely been used to extract topic metadata from different types of text data including course syllabuses (Apaza et al. 2014). However, unsupervised approaches such as LDA suffer from complex hyper-parameter tuning (Panichella 2021) and limited interpretability of *latent* KCs, creating gaps in transparency. *Wikification*, a form of entity linking (Brank, Leban, and Grobelnik 2017) has shown promise for automatically capturing the KCs covered in an educational resource (Bulathwela et al. 2020b). This technology provides *automatic, humanly-intuitive (symbolic)* representations from Wikipedia, representing *up-to-date knowledge* about *many domains*.

The possibility of recommending parts of items (contrary to an entire video or podcast) is also a fruitful research direction explored lately. From proximity-aware information retrieval (Schenkel et al. 2007) to segmenting videos to build tables of contents (Mahapatra et al. 2018), this goal has been under active research. Breaking informational videos into fragments has also shown promise in efficient previewing (Chen et al. 2018) and enabling non-linear consumption of videos (Verma et al. 2021). Recent proposals such as TrueLearn (Bulathwela et al. 2020b) demonstrate the potential of using fragment recommendation in education. Due to these reasons, we use Wikification (Brank, Leban, and Grobelnik 2017) to generate KCs that are included in each video fragment covered by the PEEKC dataset.

### Designing a Machine Learning Library

To design a user-friendly, easy-to-use, and scalable library, commonly used yet, bad design practices, such as rigidity, fragility, immobility and viscosity should be avoided (Martin 2000; Piccioni, Furia, and Meyer 2013). Many design principles are proposed to overcome these issues (Gamma et al. 1995). Many data scientists also prefer usable (adhering to known patterns), well-documented and intuitive libraries (Nadi and Sakr 2023). Besides these, designing a machine learning library entails overcoming additional challenges (e.g. data, pre-processing, models, etc.). Scikit-learn

(Pedregosa et al. 2011) proposes consistency, inspection, sensible defaults and good interface design (estimators and predictors) for building a scalable and user-friendly machine learning library. Consistency of the code interfaces significantly reduces the learning cost for users while inspection exposes relevant model parameters and public attributes to the user with easy access (Buitinck et al. 2013). The estimator interface specifies a `fit` function to provide a consistent interface to the training model and exposes the `coef_` attribute to facilitate the inspection of the internal state of the model. The predictor interface specifies the `predict` and `predict_proba` functions as methods for utilising the trained model. PyBKT, a Python-based library that implements knowledge tracing and item response theory-based learner models, also follow the same interfacing practices where function names `fit` and `evaluate` are used to train and predict (Bulut et al. 2023). Due to the time-tested and consistent design decisions that have succeeded in scikit-learn and pyBKT, we utilise the same functions to interact with the learner models in the TrueLearn library.

### Learner Modelling

Personalised learning mainly revolves around Knowledge Tracing (KT) (Bulut et al. 2023) and Item-Response Theory (IRT) (Rasch 1960) based models that use KCs in exercises to predict test success. However, these models focus on test-taking (modelling short sequences of exercise answering events) rather than consuming learning materials such as video watching. Conventional KT and IRT models do not support online learning posing scale challenges in life-long learning cases (while online counterparts exist (Bishop, Winn, and Diethe 2015)). More recently, deep-KT (Piech et al. 2015) has shown promise in superior performance. However, deep-KT models are data-hungry and lack interpretability, making them less favourable for lifelong learning, where the model needs to learn usable parameters with minimal data. Furthermore, recent studies have questioned the superior performance of deep-KT models in comparison to traditional models (Schmucker et al. 2022). Due to these reasons, we scope out batch/deep learning models and focus on data-efficient online models.

The TrueLearn family of online Bayesian learner models (Bulathwela et al. 2020b) uses implicit feedback from learners to recover their learning state. Models that capture learners' interests, knowledge, and novelty are proposed in prior work with methods to combine them as interpretable ensembles that can account for these factors simultaneously (Bulathwela, Sahan and Pérez-Ortiz, María and Yilmaz, Emine and Shawe-Taylor, John 2022). While being data efficient and privacy-preserving by design (exclusively using individual learner's interactions), TrueLearn models generate humanly intuitive learner representations inspired by Open Learner Models (OLM). This involves generating visualisations that will communicate information about learner state, promoting learner reflection by aiding learners in planning and monitoring their learning (Bull and Kay 2010). OLMs also pose challenges, since all visual presentations may not be equally understood by a wide variety of end-users. Among many visualisations used to present learner

knowledge state, user studies have shown that some visualisations are comparatively more user-friendly than others (Bull, Brusilovsky, and Guerra 2018; Bull et al. 2016). TrueLearn implements a set of tested visualisations that aid the communication and the interaction process.

Both deep KT and libraries such as pyBKT focus on predicting test-taking behaviour (Bulut et al. 2023) rather than how they would interact with an educational video. These libraries also focus on course-based learning settings where the number of KCs and learning items are limited in number. In these aspects, TrueLearn sets itself apart from the rest of the available libraries. The same reasons make TrueLearn valuable for MOOC platforms and educational video repositories that thrive to personalise videos for learning. In a world where a large number of educational videos are in circulation, we are unaware of a public, easy-to-use toolkit that can be used to incorporate educational video personalisation apart from our proposal, TrueLearn.

### Problem Setting

A learner  $\ell$  in learner population  $L$  interacting with a series of educational resources  $S_\ell \subset \{r_1, \dots, r_R\}$  where  $r_x$  are *fragments/parts* of an educational video  $v$ . The watch interactions happen over a period of  $T$  time steps,  $R$  being the total number of resources in the system. In this system with a total  $N$  unique knowledge components (KCs), resource  $r_x$  is characterised by a set of top KCs or topics  $K_{r_x} \subset \{1, \dots, N\}$ . We assume the presence  $i_{r_x}$  of KC in resource  $r_x$  and the degree  $d_{r_x} \in \{0, 1\}$  of KC coverage in the resource is observable.

The key idea is to model the probability of engagement  $e_{\ell, r_x}^t \in \{1, -1\}$  between learner  $\ell$  and resource  $r_x$  at time  $t$  as a function of the learner interest  $\theta_{\ell, \perp}^t$ , knowledge  $\theta_{\ell, \text{NK}}^t$  based on the top KCs covered  $K_{r_x}$  using their presence  $i_{r_x}$ , and depth of topic coverage  $d_{r_x}$ .

### PEEK Dataset

In this section, we describe how the Personalised Educational Engagement linked to Knowledge Components (PEEK) dataset is constructed. Figure 1 (ii) outlines the overall process of creating the PEEK dataset. PEEK uses the data from VideoLectures.Net<sup>1</sup> (VLN), a repository of scientific and educational video lectures. VLN repository records research talks and presentations from numerous academic venues (mainly AI and Computer Science). As the talks are recorded at peer-reviewed research venues, the lectures are reviewed and the material is controlled for the correctness of knowledge. Although most lectures consist of one video, some video lectures are broken into more videos (such as a long tutorial).

### Fragmenting Video Transcripts

First, the videos in VLN repository are transcribed to its native language using the *TransLectures* project<sup>2</sup>. Then, the non-English lecture videos are translated into English as we

will use English Wikipedia for entity linking. Once the transcription/translation is complete, we partition the transcript of each video into multiple *fragments* where each fragment covers approximately 5 minutes of lecture time (5000 characters). Having 5-minute fragments allows us to break the contents of a video into a more granular level while making sure that there is sufficient amounts of information while keeping fragment length at a favourable value in terms of retaining viewer engagement (Guo, Kim, and Rubin 2014).

### Wikification of Transcripts

In order to identify the Knowledge Components (KCs) that are contained in different video fragments, we use Wikification (Brank, Leban, and Grobelnik 2017). This allows annotating learning materials with humanly interpretable KCs (Wikipedia concepts) at scale with minimum human-expert intervention. This setup will make sure that recommendation strategies built on this dataset will be technologically feasible for web-scale e-learning systems.

### Knowledge Component Ranking

As per (Brank, Leban, and Grobelnik 2017), Wikification produces two statistical values per annotated KC,  $c$ , namely, *PageRank* and *Cosine Similarity* scores.

**PageRank score** is calculated by constructing a semantic graph where semantic relatedness ( $SR(c, c')$ ) between Wikipedia concept pairs  $c$  and  $c'$  in the graph are calculated using equation 1 and running PageRank on this graph.

$$SR(c, c') = \frac{\log(\max(|L_c|, |L_{c'}|) - \log(|L_c \cap L_{c'}|))}{\log |W| - \log(\min(|L_c|, |L_{c'}|))} \quad (1)$$

where  $L_c$  represents the set of Wiki concepts with inwards links to Wikipedia concept  $c$ ,  $|\cdot|$  represents the cardinality of the set and  $W$  represents the set of all Wikipedia topics. PageRank algorithm (Brin and Page 1998) leads to heavily connected Wikipedia topics (i.e. more semantically related) within the lecture to get a higher score.

**Cosine Similarity score** is used as a proxy for topic coverage within the lecture fragment (Bulathwela et al. 2020b). This score  $\cos(s_{tr}, c)$  between the *Term Frequency-Inverse Document Frequency* (TF-IDF) representations of the lecture transcript  $s_{tr}$  and the Wikipedia page  $c$  is calculated based on equation 2:

$$\cos(s_{tr}, c) = \frac{\text{TFIDF}(s_{tr}) \cdot \text{TFIDF}(c)}{\|\text{TFIDF}(s_{tr})\| \times \|\text{TFIDF}(c)\|} \quad (2)$$

where  $\text{TFIDF}(s)$  returns the TF-IDF vector of the string  $s$  while  $\|\cdot\|$  represents the norm of the TF-IDF vector.

The authors of (Brank, Leban, and Grobelnik 2017) recommend that a linearly weighted sum between the PageRank and Cosine score can be used for ranking the importance of Wikipedia concepts.

We empirically find weighting 0.8 on PageRank and 0.2 on Cosine similarity is most suitable. The ranked KCs are used to identify the five top-ranked KCs for each lecture fragment. Figure 2(ii) provides a word cloud of the most dominant KCs in the PEEK dataset. It is evident that the

<sup>1</sup>www.videolectures.net

<sup>2</sup>www.translectures.eu

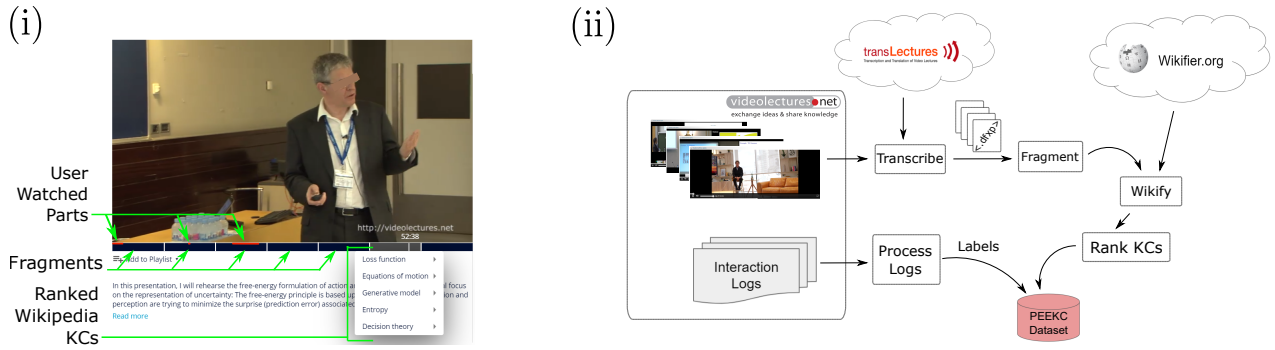


Figure 1: (i) Visual representation of the data items available in the PEEKC dataset where each video is broken into multiple, non-overlapping 5-minute fragments that are linked with ranked Wikipedia-based KCs and (ii) The flow chart presenting how the video data and the learner interaction logs from VLN repository are processed to create the PEEKC dataset.

majority of KCs associated with the lecture fragments in this dataset are related to artificial intelligence and machine learning, making this dataset ideal for training personalisation models for AI education.

### Anonymity

We restrict the final dataset to lectures with views from at least five unique users to preserve k-anonymity (Craswell et al. 2020). Also, we report the timestamp of user view events in relation to the earliest event found in the dataset obfuscating the actual timestamp. We report the smallest timestamp in the dataset  $t_0$  as 0s and any timestamp  $t_i$  after that as  $t_i - t_0$ . This allows us to publish the real order and the differences in time between events without revealing the actual timestamps. Additionally, the lecture metadata such as title and authors are not published to preserve their anonymity. The motivation here is to avoid video presenters having unanticipated effects on their reputation by associating implicit learner engagement with their content.

### Labels

The user interface of VLN website also records the video-watching behaviors of its users (see Figure 1 (i)). We create a binary target label based on *video watch time* commonly used as a proxy for video engagement in both non-educational (Covington, Adams, and Sargin 2016; Wu, Rizoiu, and Xie 2018) and educational (Guo, Kim, and Rubin 2014; Bulathwela et al. 2020b) contexts. Normalised learner watchtime  $\bar{e}_{\ell,r}^t$  of learner  $\ell$  with video fragment resource  $r_i$  at time point  $t$  is calculated as per equation 3.

$$\bar{e}_{\ell,r_i}^t = W(\ell, r_i) / D(r_i), \quad (3)$$

where  $\bar{e}_{\ell,r_i}^t \in \{0, 1\}$ ,  $W(\cdot)$  is a function returning the *watch time* of learner  $\ell$  for resource  $r_i$  and  $D(\cdot)$  is a function returning the duration of lecture fragment  $r_i$ . The final label  $e_{\ell,r_i}^t$  is derived by discretising  $\bar{e}_{\ell,r_i}^t$  where  $e_{\ell,r_i}^t = 1$  when  $\bar{e}_{\ell,r_i}^t \geq .75$  and  $e_{\ell,r_i}^t = 0$  otherwise. This rule is motivated by the hypothesis that a learner should watch approximately 4 out of 5 minutes of a video fragment in order to acquire knowledge from it (Bulathwela et al. 2020b).

Column	Description	Data Type
1	Video Lecture ID	Integer
2	Video ID	Integer
3	Part ID	Integer
4	Timestamp	Integer
5	User ID	Integer
6,8,10,12,14	Knowledge Component IDs	Integer
7,9,11,13,15	Topic Coverage	Floating Point
16	Label	Binary

Table 1: Different columns included in the PEEKC Dataset.

### Final Dataset

The final PEEKC dataset consists of 290,535 interaction events from 20,019 distinct users with at least five watch events. They engage with 8,801 unique lecture videos partitioned into 36,408 fragments (4.14 fragments per video). The learners in the dataset are divided into *Training* (14,050 learners) and *Test* (5,969 learners) datasets based on a 70:30 split. The label distribution in the dataset is also relatively balanced with only 56.35% of the labels being positive. As shown in Figure 2 (i), the majority of learners in the dataset have a relatively small number of events (under 80) making this dataset an excellent test bed for personalisation models designed to work in data-scarce environments. VLN repository mainly publishes videos relating to AI and Machine Learning, attracting a large number of learners who visit to learn about these subjects. This fact is confirmed by Figure 2 where it shows that the dataset is dominated by events with AI and ML-related KCs. The dataset is available publicly<sup>3</sup>. The set of columns in the dataset is described in table 1.

### truelearn Library

This section describes the architecture of the TrueLearn Python library. While TrueLearn provides a probability that can be mapped to a binary outcome (engaging/not engaging), the probability prediction on different videos ranks them, creating personalised recommendations.

<sup>3</sup><https://github.com/sahanbull/PEEKC-Dataset>

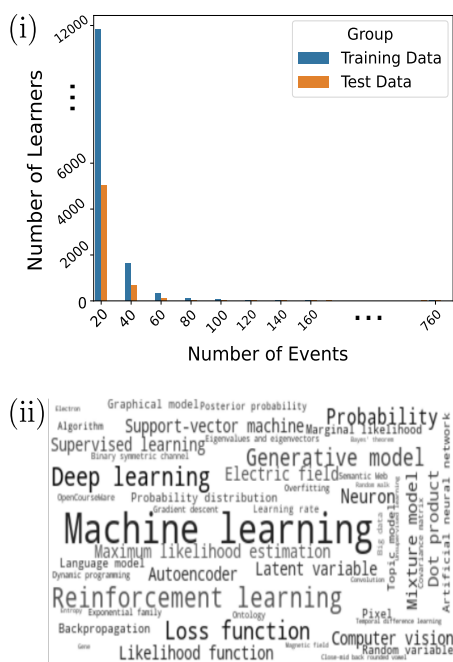


Figure 2: Characteristics of the PEEKC dataset: (i) number of learners in the training/test dataset based on the number of events in their sessions and (ii) wordcloud depicting the most frequent Wikipedia-based KCs showing the dominance of AI and ML concepts in the dataset.

## Architecture

The TrueLearn library consists of six modules.

**Datasets** The dataset module integrates tools for downloading and parsing learner engagement datasets. Currently, the PEEKC dataset is integrated.

**Pre-processing** The pre-processing module contains utility classes for extracting content representations from educational materials. The extracted representations become KCs that can be used with IRT, KT and TrueLearn models. At present, utility functions for Wikification are included.

**Models** This module houses the class that can store the learner model. In this context, the learner model refers to the data structure storing the learner state (e.g. knowledge/interest). This learner model is loosely coupled with the learning algorithms which makes this object reusable with other learning algorithms that go beyond the TrueLearn algorithms.

**Learning** This module contains machine learning algorithms that can perform training and prediction of learner engagement with transcribed videos. For training, fit function is used. For prediction, predict and predict\_proba functions are used. Currently, a set of baselines and the TrueLearn algorithms (Bulathwela, Sahana and Pérez-Ortiz, María and Yilmaz, Emine and Shawe-Taylor, John 2022) are included.

**Metrics** Classification metrics accuracy, precision, recall and F1-score are built as the task is posed as a classification task in prior work (Bulathwela et al. 2020b). The module is easily extendable to regression and ranking metrics.

**Visualisations** To effectively present the learner state, nine different visualisations shown promising in prior work on user interaction (Guesmi et al. 2022) have been developed. Figure 3 provides a preview of one of the common visualisations. Seven (out of nine) interactive visualisations allow the learner to click and hover over the output to explore more details.

## Visualising the Learner State

Our approach was guided by a thorough examination of seminal research on impactful learning visualisations. Among others, the interactive visualisations designed in our Python library takes into account the goals of self-actualisation as detailed in the EDUSS framework (Guesmi et al. 2022). Additionally, the visualisations utilise user-friendly cues and conventions (colours/intensity of colour, shape size etc.) to minimise the cognitive load. Based on user preferences found on learning visualisations (Bull et al. 2016), the i) bar plot, ii) dot plot, iii) pie plot, vi) tree plot v) radar plot, vi) rose plot vii) bubble plot viii) word plot and ix) line plot were chosen to be implemented. Figure 3 previews the learner state of one of the learners.

TrueLearn algorithms model learner skill states as Gaussian variables with mean (state estimate) and variance (estimate uncertainty). The bar plot and dot plot use the bar/dot for skill mean while mapping the uncertainty as a confidence interval. The radar plot uses the radius of the radar as the skill estimate. The pie plot and tree plot use the area of the pie to represent skill mean, while using the intensity of the colour (dark to light) for uncertainty. The rose plot, in addition to using the radius and colour intensity for mean and uncertainty respectively, uses the area of the pie to depict the number of video fragments that contributed to the skill estimate. The bubble plot and word plot use the size of the skill shape to represent the mean estimate, while the bubble plot also uses the colour intensity to depict model uncertainty (Figure 3). The line plot uses the x-axis as the time to show how a skill evolved over time.

## Experiments and Results

We use the experimental protocol used in (Bulathwela, Sahana and Pérez-Ortiz, María and Yilmaz, Emine and Shawe-Taylor, John 2022) for our experiment<sup>4</sup>.

**Baselines** We use a wide range of baselines that are i) exclusively content-based and ii) maintain a concept-based user model (Zarrinkalam et al. 2020). As content-based models, we use i) KC cosine similarity-based ( $Cosine_c$ ), Jaccard similarity based on ii) KC intersection ( $Jaccard_c$ ) and iii) user intersection ( $Jaccard_u$ ). As concept-based user models, we use iv) TF (Binary), which counts the number of times a concept was encountered and v) TF (Cosine),

<sup>4</sup>Further experiment details and other supplementary information found at: <https://arxiv.org/abs/2401.05424>



Model	Acc.	Prec.	Rec.	F1
Cosine	55.08	57.86	58.45	54.06
Jaccard <sub>c</sub>	55.46	57.81	60.36	55.03
Jaccard <sub>u</sub>	64.06	57.85	72.76	61.22
TF (Binary)	55.19	56.71	66.60	57.38
TF (Cosine)	55.11	56.75	65.95	57.11
KT	54.99	53.25	28.56	34.51
TrueLearn Interest	58.13	52.08	<i>78.61*</i>	63.00*
TrueLearn Novelty	<i>64.78*</i>	<i>58.52*</i>	<b>80.91*</b>	<b>65.53*</b>
TrueLearn INK	<b>78.32*</b>	<b>64.32*</b>	64.03	<i>64.00*</i>

Table 2: Performance of TrueLearn algorithms is evaluated against the baselines using Precision (Prec.), Recall (Rec.) and F1 Score (F1). The best and second best performance is indicated in bold and italic faces respectively. TrueLearn models outperform the best baseline most times ( $p < 0.01$  in a one-tailed paired t-test are marked with \*).

which aggregates the cosine scores for skills in PEEKC dataset over time and vi) online Knowledge Tracing model (KT) (Bishop, Winn, and Diethe 2015).

**TrueLearn Models** We use the three TrueLearn models implemented in the library, namely, i) TrueLearn Interest, capturing interests, ii) TrueLearn Novelty, capturing knowledge and novelty and iii) TrueLearn INK, combining interests, knowledge and novelty states (Bulathwela et al. 2020a).

**Data and Evaluation** For each learner, their engagement at time  $t$  is predicted using its events at times 1 to  $t - 1$ . Hold-out validation (70% train/ 30% test) technique is used in our experiment where the training data in PEEKC is used for hyperparameter tuning. The best hyperparameter combination based on the F1-Score is identified and used with the test set to evaluate the reported performance. Since the engagement is labelled as a binary label in the PEEKC dataset, accuracy, precision, recall, and F1 score are reported.

### Empirical Evaluation

The results are reported in Table 2. Our experiments i) guarantee the correctness of the library implementation and ii) demonstrate the predictive capabilities of the web-scale online learning models with comparable baselines.

### Discussion

The contributions in this work span over a novel dataset, an OLM library and empirical experiments.

### PEEKC for AI Education

As per the wordcloud in Figure 2 (ii), the videos in the PEEKC dataset are swamped with AI and ML-related KCs. PEEKC’s data source, VLN repository extensively visits AI conferences which causes this effect. But, this makes PEEKC a perfect dataset to understand in-the-wild video-based knowledge acquisition, making this dataset an excellent resource for training personalisation models for AI education as we have done in table 2. While we haven’t demonstrated here, the dataset is potentially valuable for unsupervised tasks such as pre-requisite identification and hierarchical structuring of concepts in AI. A key benefit of PEEKC

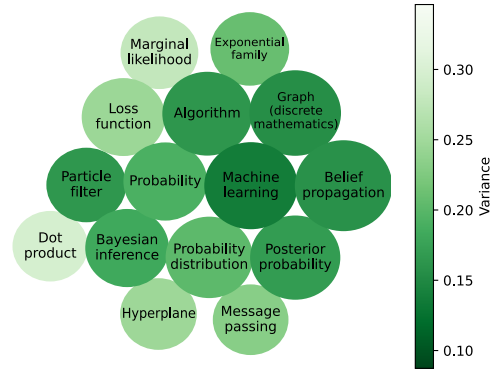


Figure 3: The 15 most knowledge acquired KCs of a learner in a bubble plot. The size of the circle aligns with the KC mean and the intensity of the colour maps to the variance.

dataset is the ability to use the methodology in Figure 1 with VLN repository to create larger datasets in the future.

### TrueLearn Performance and Visualisations

Table 2 clearly shows the superiority of the TrueLearn model implementations in comparison to the comparative baselines aligning with prior work’s evidence (Bulathwela, Sahan and Pérez-Ortiz, María and Yilmaz, Emine and Shawe-Taylor, John 2022). The most expensive, `fit` function measures  $0.005 \pm 0.0036s$  to execute in an Apple M1 3.2GHz CPU. TrueLearn states of each learner in the dataset can be constructed in parallel because there are no inter-learner dependencies, also making the model privacy-preserving by design. The performance in table 2 coupled with event counts in Figure 2 shows the data efficiency of TrueLearn, which is able to learn from a very small number of events.

The library is designed to enable learners to effortlessly generate dynamic and static visualisations, promoting a learner-centric, self-regulated study experience (Bull and Kay 2008). Figure 3 previews the learner knowledge state from one of the learners in the PEEKC test data. It is seen that the visualisation demonstrates how the user is acquiring knowledge in AI-related topics, specifically, in the area of Bayesian modelling. The EDUSS framework-inspired visualisations support development by translating skill-level predictions into visual progress indicators encouraging continuous learning. Understanding is enhanced by providing insights into learners’ interaction patterns, promoting self-awareness (Guesmi et al. 2022). The transparency of the user model allows learners to scrutinise their learning progress and engage critically with the underlying data (while ways to provide this feedback to the model is still an open research area). Finally, by making the visualizations shareable, the library fosters social interaction, adding a community dimension to the learning experience and improving its usability to other stakeholders of education (e.g. parents and teachers).

## Library Design, Stability and Maintainability

The adherence to prior work-inspired good practices and API design makes using the TrueLearn library developer-friendly for both offline experimentation and e-learning system integration. Using a collection of base classes that define a common interface and shared functionality similar to scikit-learn-like estimators with `fit` and `predict` functions (Buitinck et al. 2013) reduces the learning curve while the design allows elaborate type and value checks when the hyperparameters of the classifiers are modified, ensuring the robustness of the classifier implementation. The decoupling of the feature-engineering from modelling modules allows users to extend the capability of the library with new KC extraction functions (beyond Wikification), new online learner models and open learner visualisations without having to worry about interdependencies. Detailed instructions are provided to developers with style guides and design principals<sup>5</sup> to ease contribution. The core research team is committed to maintaining the library in the future while providing further guidance and code reviews when extending the capabilities of the library. Furthermore, TrueLearn benefits from 100% test coverage achieved through a combination of integration, unit, and documentation tests.

Code consistency and readability are further enhanced by following the PEP 8 guidelines (van Rossum, Warsaw, and Coghlan 2001), which define a set of best practices for Python code. Extensive documentation of the modules and classes with in-context code examples describes relevant information for both a potential developer and a contributor to familiarise themselves with the library. Developers have already started adapting this library to learning platforms (Bulathwela, Kreitmayer, and Pérez-Ortiz 2020). We aim to objectively assess their experience by surveying them (Piccioni, Furia, and Meyer 2013; Nadi and Sakr 2023).

## Relevance, Impact and Limitations

Modelling learner state in a humanly intuitive manner, requiring minimal data and exclusively relying on individual user actions, TrueLearn offers a transparent learner model that respects the privacy of its users and can scale to lifelong education. The development of the TrueLearn library aims to provide both the research and developer communities with the opportunity to seamlessly use the TrueLearn family of models in their work. The learner models utilise Wikipedia-based entity linking to create KCs based on a publicly available knowledge base. The content annotation can also scale to thousands of materials created in different modalities (video, text, audio etc.).

The impact of TrueLearn is two-fold. For developers and researchers, the TrueLearn library employs a design that conforms with popular machine learning libraries. The documentation is extensive and contains detailed examples that help the implementation. For developers and educators, probabilistic graphical models that are data efficient and humanly intuitive are available to be used in their downstream systems. The engagement predictions (between values 0 and 1) can be used to rank videos for personalised learning.

<sup>5</sup>Contributing: <https://truelearn.readthedocs.io/en/latest/dev/>

Combined with the PEEKC dataset, hyperparameters for a new system can be trained beforehand and deployed in a new online learning platform. The online learning algorithm updates the learner state in real-time helping better personalisation. A platform implementing TrueLearn can scale to a large population of users and support them through lifelong education due to the large number of KCs it can support.

While getting inspired by scikit-learn library, the learning algorithms in TrueLearn library are not compatible with some helper functions available in scikit-learn (such as grid search) and pandas libraries at this point. Building seamless compatibility with these utilities will enable the TrueLearn library to be adopted by a wider audience while minimising the development effort required to support such powerful features. While the visualisations implemented are time-tested (Bull et al. 2016; Guesmi et al. 2022), their success with TrueLearn representations has room for rigorous understanding via user studies. The exclusive support of online learning algorithms can also be seen as a limitation of the current library as many batch learning algorithms are proposed for educational recommendation and engagement modelling (Lin and Chi 2016; Pardos et al. 2012). While learner engagement is a prerequisite for learning, it is noteworthy that learner engagement doesn't imply learning. The library also does not support state-of-the-art deep learning algorithms (Piech et al. 2015; Pardos and Jiang 2020) that may be useful where interpretability and learner state visualisation are not the top priority.

## Conclusion

This work presents *PEEK*C dataset and *TrueLearn* Python library, creating a valuable toolbox for engagement modelling with AI-related educational videos. The library contains several online learning models, which model multiple factors influencing learner engagement. It also packages a set of visualisations that can be used to interpret the learner's interest/knowledge state. The learner representations and state visualisations are comparable to outputs of knowledge tracing models, except TrueLearn uses watch time interactions rather than relying on test taking. The empirical results demonstrate that the implementation of the library achieves similar performance to the prior work. The new implementation encourages educational data mining practitioners to use this library to incorporate educational video recommendations in e-learning systems. Researchers are encouraged to extend this library with new datasets and online learning algorithms for learner engagement modelling.

The immediate future work entails improving learner state visualisations via user studies. Integrating the library into a real-world e-learning platform (Perez-Ortiz et al. 2021) is a top priority. Extending the current framework to podcasts and other information content while incorporating other feedback forms like educational questions (Bulathwela, Muse, and Yilmaz 2023; Fawzi, Amini, and Bulathwela 2023) remains in the future roadmap. In the long term, we aim to add more general informational recommendation algorithms to the library and mobilise the research community to contribute various models, pre-processing techniques and evaluation metrics that the library can benefit from.

## Acknowledgments

This work is partially supported by the European Commission-funded project "Humane AI: Toward AI Systems That Augment and Empower Humans by Understanding Us, our Society and the World Around Us" (grant 952026), the X5GON project funded from the EU's Horizon 2020 research programme grant No 761758 and EU Erasmus+ project (Encore+) 621586-EPP-1-2020-1-NO-EPPKA2-KA. An extended version of this paper with more details about the dataset and the experiments can be found at <https://arxiv.org/abs/2401.05424>

## References

- Apaza, R. G.; Cervantes, E. V. V.; Quispe, L. V. C.; and Luna, J. E. O. 2014. Online Courses Recommendation based on LDA. In *Symposium on Information Management and Big Data*.
- Bishop, C.; Winn, J.; and Diethe, T. 2015. *Model-Based Machine Learning*. Early access version (<http://www.mbmlbook.com/>). Accessed: 2019-05-23.
- Bloom, B. S. 1984. The 2 Sigma Problem: The Search for Methods of Group Instruction as Effective as One-to-One Tutoring. *Educational Researcher*, 13(6): 4–16.
- Brank, J.; Leban, G.; and Grobelnik, M. 2017. Annotating Documents with Relevant Wikipedia Concepts. In *Proc. of Slovenian KDD Conf. on Data Mining and Data Warehouses (SiKDD)*.
- Brin, S.; and Page, L. 1998. The anatomy of a large-scale hypertextual Web search engine. In *Proc. of Int. Conf. on World Wide Web*.
- Buitinck, L.; Louppe, G.; Blondel, M.; Pedregosa, F.; Mueller, A.; Grisel, O.; Niculae, V.; Prettenhofer, P.; Gramfort, A.; Grobler, J.; Layton, R.; VanderPlas, J.; Joly, A.; Holt, B.; and Varoquaux, G. 2013. API design for machine learning software: experiences from the scikit-learn project. *CoRR*, abs/1309.0238.
- Bulathwela, S.; Kreitmayer, S.; and Pérez-Ortiz, M. 2020. What's in It for Me? Augmenting Recommended Learning Resources with Navigable Annotations. In *Proceedings of the 25th International Conference on Intelligent User Interfaces Companion, IUI 20*.
- Bulathwela, S.; Muse, H.; and Yilmaz, E. 2023. Scalable Educational Question Generation with Pre-trained Language Models. In *International Conference on Artificial Intelligence in Education*, 327–339. Springer.
- Bulathwela, S.; Pérez-Ortiz, M.; Mehrotra, R.; Orlic, D.; de la Higuera, C.; Shawe-Taylor, J.; and Yilmaz, E. 2020a. Report on the WSDM 2020 Workshop on State-based User Modelling (SUM'20). volume 54. ACM.
- Bulathwela, S.; Pérez-Ortiz, M.; Yilmaz, E.; and Shawe-Taylor, J. 2020b. TrueLearn: A Family of Bayesian Algorithms to Match Lifelong Learners to Open Educational Resources. In *AAAI Conference on Artificial Intelligence, AAAI 20*.
- Bulathwela, S.; Pérez-Ortiz, M.; Holloway, C.; Cukurova, M.; and Shawe-Taylor, J. 2024. Artificial Intelligence Alone Will Not Democratise Education: On Educational Inequality, Techno-Solutionism and Inclusive Tools. *Sustainability*, 16(2).
- Bulathwela, Sahan and Pérez-Ortiz, María and Yilmaz, Emine and Shawe-Taylor, John. 2022. Power to the Learner: Towards Human-Intuitive and Integrative Recommendations with Open Educational Resources. *Sustainability*, 14(18).
- Bull, S.; Brusilovsky, P.; and Guerra, J. 2018. Which Learning Visualisations to Offer Students? In Pammer-Schindler, V.; Pérez-Sanagustín, M.; Drachsler, H.; Elferink, R.; and Scheffel, M., eds., *Lifelong Technology-Enhanced Learning*, 524–530. Cham: Springer International Publishing. ISBN 978-3-319-98572-5.
- Bull, S.; Brusilovsky, P.; Guerra, J.; and Araújo, R. 2016. Individual and Peer Comparison Open Learner Model Visualisations to Identify What to Work On Next. In *Late-breaking Results, Posters, Demos, Doctoral Consortium and Workshops Proceedings of the 24th ACM Conference on User Modeling, Adaptation and Personalisation (UMAP 2016)*.
- Bull, S.; and Kay, J. 2008. Metacognition and open learner models. In *The 3rd workshop on meta-cognition and self-regulated learning in educational technologies, at ITS2008*, 7–20.
- Bull, S.; and Kay, J. 2010. *Open Learner Models*, 301–322. Springer Berlin Heidelberg. ISBN 9783642143632.
- Bulut, O.; Shin, J.; Yildirim-Erbasli, S. N.; Gorgun, G.; and Pardos, Z. A. 2023. An Introduction to Bayesian Knowledge Tracing with pyBKT. *Psych*, 5(3): 770–786.
- Chen, J.; Chen, X.; Ma, L.; Jie, Z.; and Chua, T.-S. 2018. Temporally grounding natural sentence in video. In *Proceedings of the 2018 conference on empirical methods in natural language processing*, 162–171.
- Choi, Y.; Lee, Y.; Shin, D.; Cho, J.; Park, S.; Lee, S.; Baek, J.; Bae, C.; Kim, B.; and Heo, J. 2020. Ednet: A large-scale hierarchical dataset in education. In *International Conference on Artificial Intelligence in Education*, 69–73. Springer.
- Corbett, A. T.; and Anderson, J. R. 1994. Knowledge tracing: Modeling the acquisition of procedural knowledge. *User modeling and user-adapted interaction*, 4(4): 253–278.
- Covington, P.; Adams, J.; and Sargin, E. 2016. Deep Neural Networks for YouTube Recommendations. In *Proc. of ACM Conf. on Recommender Systems*.
- Craswell, N.; Campos, D.; Mitra, B.; Yilmaz, E.; and Billerbeck, B. 2020. ORCAS: 20 Million Clicked Query-Document Pairs for Analyzing Search. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management, CIKM '20*, 2983–2989. New York, NY, USA: Association for Computing Machinery.
- Deng, W.; Zhu, P.; Chen, H.; Yuan, T.; and Wu, J. 2023. Knowledge-aware sequence modelling with deep learning for online course recommendation. *Information Processing & Management*, 60(4): 103377.
- Fawzi, F.; Amini, S.; and Bulathwela, S. 2023. Small Generative Language Models for Educational Question Generation. In *In Proc. of the NeurIPS Workshop on Generative AI for Education (GAIED)*.



- Gamma, E.; Helm, R.; Johnson, R.; and Vlissides, J. M. 1995. *Design patterns: elements of reusable object-oriented software*. Pearson Deutschland GmbH.
- Guesmi, M.; Chatti, M. A.; Tayyar, A.; Ain, Q. U.; and Joarder, S. 2022. Interactive Visualizations of Transparent User Models for Self-Actualization: A Human-Centered Design Approach. *Multimodal Technologies and Interaction*, 6(6).
- Guo, P. J.; Kim, J.; and Rubin, R. 2014. How Video Production Affects Student Engagement: An Empirical Study of MOOC Videos. In *Proc. of the First ACM Conf. on Learning @ Scale*.
- Lin, C.; and Chi, M. 2016. Intervention-BKT: Incorporating Instructional Interventions into Bayesian Knowledge Tracing. In Micarelli, A.; Stamper, J.; and Panourgia, K., eds., *Proc. of Int. Conf. on Intelligent Tutoring Systems*.
- Mahapatra, D.; Mariappan, R.; Rajan, V.; Yadav, K.; A., S.; and Roy, S. 2018. VideoKen: Automatic Video Summarization and Course Curation to Support Learning. In *Companion Proceedings of the The Web Conference 2018, WWW '18*, 239–242. Republic and Canton of Geneva, CHE: International World Wide Web Conferences Steering Committee. ISBN 9781450356404.
- Martin, R. C. 2000. Design principles and design patterns. *Object Mentor*, 1(34): 597.
- Nadi, S.; and Sakr, N. 2023. Selecting third-party libraries: the data scientist's perspective. *Empirical Software Engineering*, 28(1): 15.
- Panichella, A. 2021. A Systematic Comparison of search-Based approaches for LDA hyperparameter tuning. *Information and Software Technology*, 130: 106411.
- Pardos, Z. A.; Gowda, S. M.; Baker, R. S.; and Heffernan, N. T. 2012. The Sum is Greater than the Parts: Ensembling Models of Student Knowledge in Educational Software. *SIGKDD Explor. Newsl.*, 13(2): 37–44.
- Pardos, Z. A.; and Jiang, W. 2020. Designing for Serendipity in a University Course Recommendation System. In *Proceedings of the Tenth International Conference on Learning Analytics & Knowledge, LAK '20*, 350–359. New York, NY, USA: Association for Computing Machinery. ISBN 9781450377126.
- Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; et al. 2011. Scikit-learn: Machine learning in Python. *the Journal of machine Learning research*, 12: 2825–2830.
- Perez-Ortiz, M.; Dormann, C.; Rogers, Y.; Bulathwela, S.; Kreitmayer, S.; Yilmaz, E.; Noss, R.; and Shawe-Taylor, J. 2021. X5Learn: A Personalised Learning Companion at the Intersection of AI and HCI. In *26th International Conference on Intelligent User Interfaces - Companion, IUI '21 Companion*, 70–74. New York, NY, USA: Association for Computing Machinery. ISBN 9781450380188.
- Piccioni, M.; Furia, C. A.; and Meyer, B. 2013. An empirical study of API usability. In *2013 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*, 5–14. IEEE.
- Piech, C.; Bassen, J.; Huang, J.; Ganguli, S.; Sahami, M.; Guibas, L. J.; and Sohl-Dickstein, J. 2015. Deep Knowledge Tracing. In *Advances in Neural Information Processing Systems*.
- Rasch, G. 1960. *Probabilistic Models for Some Intelligence and Attainment Tests*, volume 1.
- Schenkel, R.; Broschart, A.; Hwang, S.; Theobald, M.; and Weikum, G. 2007. Efficient text proximity search. In *International Symposium on String Processing and Information Retrieval*, 287–299. Springer.
- Schmucker, R.; Wang, J.; Hu, S.; and Mitchell, T. 2022. Assessing the Performance of Online Students - New Data, New Approaches, Improved Accuracy. *Journal of Educational Data Mining*, 14(1): 1–45.
- Selent, D.; Patikorn, T.; and Heffernan, N. 2016. ASSISTments Dataset from Multiple Randomized Controlled Experiments. In *Proc. of the Third (2016) ACM Conf. on Learning @ Scale, L@S '16*. New York, NY, USA: Association for Computing Machinery. ISBN 9781450337267.
- van Rossum, G.; Warsaw, B.; and Coghlan, N. 2001. Style Guide for Python Code. PEP 8.
- Verma, G.; Nalamada, T.; Harpavat, K.; Goel, P.; Mishra, A.; and Srinivasan, B. V. 2021. Non-Linear Consumption of Videos Using a Sequence of Personalized Multimodal Fragments. In *26th International Conference on Intelligent User Interfaces, IUI '21*, 249–259. New York, NY, USA: Association for Computing Machinery. ISBN 9781450380171.
- Wang, Z.; Lamb, A.; Saveliev, E.; Cameron, P.; Zaykov, Y.; Hernandez Lobato, J. M.; Turner, R.; Baraniuk, R.; Barton, C.; Peyton Jones, S.; Woodhead, S.; and Zhang, C. 2020. Diagnostic Questions: The NeurIPS 2020 Education Challenge.
- Wang, Z.; Tschitschek, S.; Woodhead, S.; Hernández-Lobato, J. M.; Peyton, J. S.; Baraniuk, R. G.; and Zhang, C. 2021. Educational Question Mining At Scale: Prediction, Analysis and Personalization. In *Symposium on Educational Advances in Artificial Intelligence (AAAI-EAAI)*.
- Wu, S.; Rizozi, M.; and Xie, L. 2018. Beyond Views: Measuring and Predicting Engagement in Online Videos. In *Proc. of the Twelfth Int. Conf. on Web and Social Media*.
- Yu, J.; Luo, G.; Xiao, T.; Zhong, Q.; Wang, Y.; Feng, W.; Luo, J.; Wang, C.; Hou, L.; Li, J.; et al. 2020a. MOOC-Cube: a large-scale data repository for NLP applications in MOOCs. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 3135–3142.
- Yu, Y.; Karlgren, J.; Bonab, H.; Clifton, A.; Tanveer, M. I.; and Jones, R. 2020b. Spotify at the TREC 2020 Podcasts Track: Segment Retrieval. In *Proceedings of the Twenty-Ninth Text REtrieval Conference (TREC 2020)*.
- Zarrinkalam, F.; Faralli, S.; Piao, G.; and Bagheri, E. 2020. *Extracting, mining and predicting users' interests from social media*. Hanover, MD: Now Foundations and Trends.