

A Data-Driven Approach for Inferring Student Proficiency from Game Activity Logs

Mohammad H. Falakmasir^{*,#} José P. González-Brenes^{*} Geoffrey J. Gordon[§] Kristen E. DiCerbo^{*}
^{*}School Research [#]Intelligent Systems Program [§]Machine Learning
Pearson University of Pittsburgh Carnegie Mellon University
{jose.gonzalez-brenes, falakmasir@pitt.edu ggordon@cs.cmu.edu
kristen.dicerbo}@pearson.com

ABSTRACT

Traditional assessments are important in education because they allow collecting evidence about student progress. Unfortunately, they can be very tedious to their stakeholders. In contrast, invisible assessment unobtrusively gathers performance data from students' daily activities to make inference about student relevant competency. We present a novel data analysis pipeline, Student Proficiency Inferer from Game data (SPRING), that allows modeling game playing behavior in educational games. Unlike prior work, SPRING is a fully data-driven method that does not require costly domain knowledge engineering. We validate our method using data collected from students playing 11 educational mini-games. Our results suggest that SPRING is accurate to predict Math assessments ($R^2 = 0.55$, Spearman $\rho = 0.82$).

ACM Classification Keywords

K.3.1 Computer Uses in Education.

Author Keywords

Educational Games, Student Modeling, Stealth Assessment

INTRODUCTION

Educational assessments are important because they collect evidence about whether the teaching goals are being met. Unfortunately, the process of administering assessments is usually disconnected to the learning environment, and it is often disruptive to the classroom. In many developed countries, students now find themselves spending increasing amounts of time preparing and taking tests instead of learning [7]. For example, a survey of the current state of testing in America revealed that students are taking an average of 113 standardized tests between pre-K and highschool [8]. For these reasons, it is not surprising that the recent political climate and the general population have been weighing in on the question of whether students are being tested too much [8].

According to Evidence Centered Design (ECD) [13], the goal of assessment is to characterize the strength of evidence regarding claims one wants to make about individuals or groups. Therefore, the assessment process involves identifying, organizing, or creating activities for students so we may observe that evidence. An interesting alternative to traditional summative assessment is invisible (or stealth) assessment [15], where the evidence is gathered from learners unobtrusively from the digital interactions of their ongoing activities. This data is used to understand claims regarding what students know and what they can do [16]. Stealth assessment is also intended to reduce or eliminate test anxiety, while not sacrificing validity and reliability [14].

A promising opportunity for invisible assessment is using log data collected from educational games. Unfortunately, engineering a system that parses logs is costly and time-consuming. For example, prior work [15, 16] has relied extensively on subject matter expertise to define a student model in form of a Bayesian Network to build both the competency model, and to extract features of the performance data.

Our motivation is that invisible assessment from game data may become more accurate and cheaper to implement if the domain knowledge engineering could be automated by a data-driven process. Game data is often logged with what we call a *slot and filler* structure. In Table 1, we show a simplified example of a real educational game log that uses slot and filler structure. The slots are discrete sets of events that are initiated by the learner. Each slot may accept zero to multiple fillers. Each filler represents a value of a property of the slot event. For example, a “move” event that represents the learner moved an object in the screen, may have an *x* and *y* fillers, to represent the target position in two-dimensional space.

Id	User Id	Event Name	Event Data
1	ABC	Game Start	{ }
2	ABC	Move Object	{X:363,Y:82}
3	ABC	Move Object	{X:361,Y:54}
4	ABC	Open Toolbox	{ }
4	ABC	Activate Tool	{tool:gluumi}
5	ABC	Use Gluumi	{sizeGluedTo:8,sizeNew:9}
...			

Table 1. An example fragment of a log from an educational game.

I renamed X and Y for simplicity, hope that's ok

Paste the appropriate copyright statement here. ACM now supports three different copyright statements:

- ACM copyright: ACM holds the copyright on the work. This is the historical approach.

- License: The author(s) retain copyright, but ACM receives an exclusive publication license.

- Open Access: The author(s) wish to pay for the work to be open access. The additional fee must be paid to ACM.

This text field is large enough to hold the appropriate release statement assuming it is single spaced.

Every submission will be assigned their own unique DOI string to be included here.

Conventional machine learning algorithms cannot input slot and filler data, as they work in structures called *feature vectors* or *sequences*. Prior work [4] provides a review of the difference between feature vectors and sequences. A feature vector representation requires mapping an observation onto a fixed number of features. It is not trivial to map sequences of student actions that can be of an arbitrary length into a feature vector that needs to be of a predetermined dimension. Traditional feature vector classifiers, like logistic regression or decision trees, that are used in off-the-shelf data mining packages, such as Weka [5], cannot use slot and filler data as input. In contrast, machine learning algorithm that allow sequential representations, like Hidden Markov Models (HMM), require a parametric model with the same number of dimensions for every observation. This does not occur in slot and filler structures. For example, in Table 1, the Move Object slot requires X and Y fillers, while the Use Gluumi slot requires size fillers.

In this paper, we propose *Student P*rofficiency *I*nferrer from *G*ame data (SPRING), a novel data analysis pipeline that models game playing behavior. SPRING allows modeling raw data in slot and filler structure. We demonstrate our approach on real student game playing data. Our experiments suggest that SPRING can be used to predict student proficiency without costly domain expertise.

SPRING FRAMEWORK

SPRING is designed in a way that can capture sequential decision making process of students in a way that is representative of their mastery with minimum reliance on expert knowledge. Our data analysis pipeline receives raw data of student interactions with the different levels of educational game in slot and filler format along with their post-test results and creates a regression model for predicts the post-test score. Algorithm 1 describes the three main steps of our pipeline: discretization, sequence modeling, and regression. **Explain briefly in words what the whole thing does.**

We now explain the different steps in detail.

Procedure 1 The SPRING algorithm

Input: : A log file $\mathbf{L}_{g,s}$ of slot-filler structure for each game g and student s , a test y_s for each student

1:	$\langle \mathbf{x}_{g,s}, c_s \rangle \leftarrow \text{Cluster_students}(\mathbf{L}_{g,s}, y_s)$	} Discretization
2:	$\mathbf{x}_{g,s} \leftarrow \text{Discretize}(\mathbf{L}_{g,s})$	
3:	for each game g do	
4:	for each cluster c_s do	} Sequence Modeling
5:	$\theta_{g,c_s} \leftarrow \text{Learn_HMM}(\mathbf{x}_{g,s})$	
6:	end for	
7:	end for	} Regression
8:	Prediction y_s from $x_{g,s}, \theta_{g,c_s}$	

Discretization

Each section should refer to line numbers in the algorithm

Make sure the notation is consistent between algorithm and steps

This step of the data analysis pipeline receives time-stamped student interactions as input and generates a sequence of multinomial observations for each student in each game

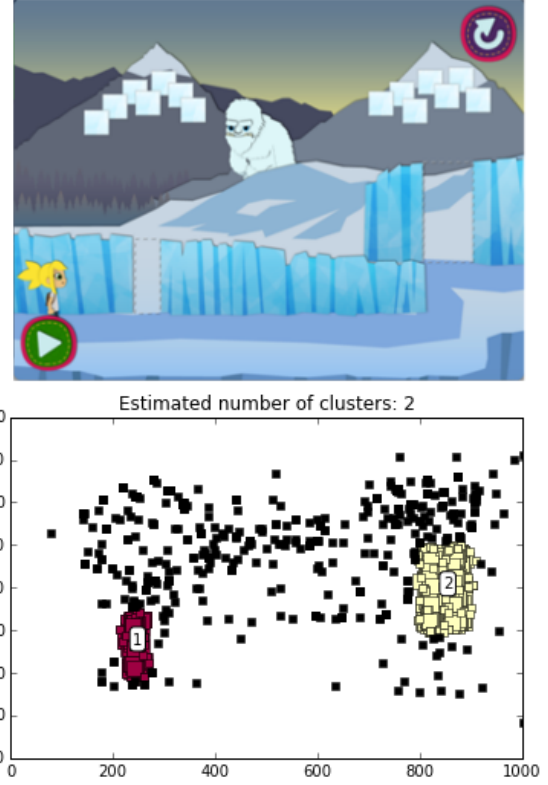


Figure 1. The clusters detected for game level 2, *None Shall Pass!*, using DBSCAN method. We transform each movement action into corresponding cluster id during the discretization step.

level. In order to transform the slot and filler observations into multinomial, we use DBSCAN clustering algorithm [2] on the filler dimensions. DBSCAN was a good candidate in our implementation because student interaction with our educational game was dominated by moving objects in a 2 dimensional space. We tried different unsupervised and non-parametric methods, however, we ended up using DBSCAN because it can find arbitrarily shaped clusters and does not require one to specify the number of clusters in the data a-priori. It is also robust to outliers which in our case can be interpreted as noisy movements that are not representative of learning competency. We considered outliers as a separate cluster.

Figure 1 shows the clusters detected for one of the game levels. In this game level (screen shot on top of Figure 1), students had to move the ice cubes from the mountain tops onto the designated area. As it is illustrated in Figure 1 bottom, DBSCAN successfully identifies the designated regions as clusters and movement outside of them as outliers.

After running DBSCAN on the filler dimensions, we trained a K-Nearest Neighbor Classifier on the cluster labels in order to convert the sequence of slot and filler observations into a two dimensional array of multinomial observations. We also used this classifier to transform any unseen sequence of

interactions for our empirical evaluation phase. The details of the discretization phase is demonstrated in Procedure 2.

Procedure 2 The Discretization Step of SPRING

```

procedure DISCRETIZE(S)
  A = empty dictionary of slots and possible fillers
  for each sequence s in S do
    for each action a in s do
      if a.filler  $\neq \emptyset$  then
        A[a.slot].append(a.filler)
      end if
    end for
  end for

  for each slot and fillers tuple in A do
    clustersIDs = DBSCAN(a.fillers)
    C[a.slot] = K-NN classifier on fillers and clusterIDs
  end for

  D =  $\square\square$   $-2D$  array of multinomial observations
  for each sequence s in S do
    P[i] = s.post-test
    for each action a in s do
      D[i,j] = C[a.slot].predict(filler)
    end for
  end for
  return D
end procedure

```

Sequence Modeling

In the modeling phase, we hypothesized that high- and low-performing students have similar usage patterns that are representative of their sequential decision-making process. Our model should be able to successfully capture this difference across different game levels and use it to predict post-test scores. Hidden Markov Models are good candidate for the task of unsupervised analysis of sequential data. Given the sequence of student actions, we aim to infer meaningful sequence of (latent) states, which describe the process that generated the actions, along with statistical patterns that can describe and distinguish those states. Since a large portion of student interactions are based on the logic of the game, we had no idea about the number of states and what they mean. As a result, we used the Hierarchical Dirichlet Process HMM (HDP-HMM) [3], which allows state spaces of unknown size to be learned from data. HDP-HMM defines a *hierarchical Dirichlet process* prior on transition matrices over countably infinite state spaces and is able to make a principled choice of how many states it needs based on the complexity of its training data. For details on training methods please refer to [3].

We used the output of the discretization step (2D array of multinomial observations) and trained two HDP-HMMs, one for high-performing and the other for low-performing students in each game level. The two models can be considered as a stochastic representation of the sequence of actions and we can use them to infer the likelihood of any arbitrary sequence as a feature for the regression step.

For simplicity, here we assume only two clusters for students being in ei

Talk about Forward backward algorithm, and how you calculate theta!!!

Regression

For each student s in each game level g , we calculate the difference ($d_{s,g}$) of the likelihood of belonging to high performing group minus the likelihood of belonging to the low performing group. We calculate this likelihood by estimating the Forward-Backward probabilities on each student sequence of actions based on the two HMMs parameter we estimated (whether the student is in the high performing group or in the low performing group):

$$d_{s,g} = \theta_{g,\text{high}} - \theta_{g,\text{low}} \quad (1)$$

We use a linear regression model for predicting the post-test scores:

$$\hat{y}_s(\beta) = \sum_g \beta_g \cdot d_{s,g} + \beta_0 \quad (2)$$

Here, β_0 is just an intercept for the model. We can optimize the parameters using **what algorithm?**. We experimented with different regularization methods, but only report **LASSO ??** **cite** as it worked best in our preliminary comparisons:

$$\beta^* = \underset{\beta}{\operatorname{argmin}} ||y_s - \hat{y}_s(\beta)||_2 + \lambda \cdot ||\beta||_1 \quad (3)$$

EMPIRICAL EVALUATION

Game Environment

Alice in AreaLand is an educational game developed for research purposes. It focuses on teaching and assessing geometric measurement, specifically the understanding of area, among 6th grade students. The game targets three main stages in the development of area: 1) area unit iteration, 2) use of unit squares to measure area, and 3) use of composites to measure area. The current version has 11 game levels. A simple student scenario involves covering a 2D area with smaller unit squares placed end-to-end in non-overlapping fashion, combining the single squares into rows or columns, and then determining the number of rows or columns needed. Figure 2 shows a screenshot of one game level.

Throughout the game, *Alice* is accompanied by *Flat Cat* – an assistant character who provides feedback and scaffolding to the player in the beginning of each game level and upon request when students push a hint button (represented by two magnifiers at the bottom center of Figure 2). Earlier game levels are designed for students to learn about area unit iteration and usually require them to cover a number of predefined areas with unit squares (not necessarily in a non-overlapping fashion). By advancing through game levels, students are presented with three tools: *Gluumi* for combining unit squares by gluing them together; *Multy*, for making copies of different objects; and *Esploda* for breaking compound shapes into single units. There is no limit for completing a game level regarding time or number of actions students may execute. The students press the *Go Alice* button (bottom left corner of Figure 2) if

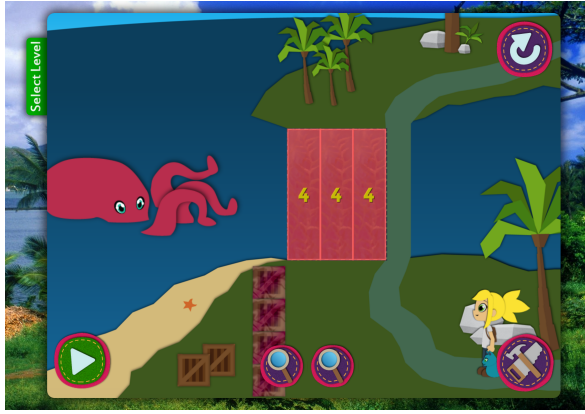


Figure 2. A screenshot of hint provided in game level 11, *You Kraken Me Up!*, in *Alice in AreaLand*. Students should combine four squares into a column and create three copies of the column to cover the designated area and prevent the octopus from attacking *Alice* while she crosses the bridge.

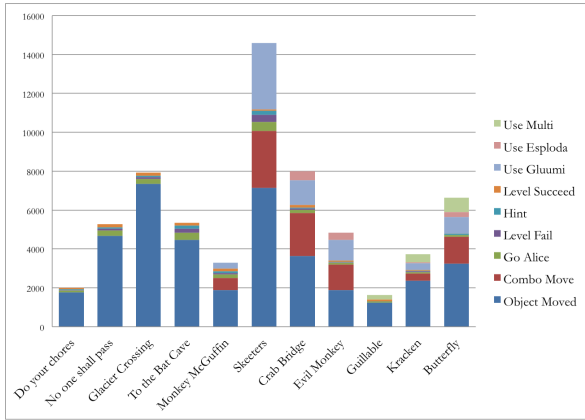


Figure 3. Frequency of Events in Each Game Level

they deem their performance to be satisfactory for *Alice* to proceed. Based on the covered area and the arrangement of the tiles, they either advance to the next level or receive a feedback and stay in the same level.

Dataset

Our dataset consists of time-stamped interactions of 129 students in 11 game levels. For 77 students, we also have post-test scores from a paper-based exam with 20 questions in the 3 skills of geometric measurement. In total, there are 88,458 events recorded in the dataset from 1,510 game sessions, meaning that student tried some of the game levels for multiple times. Based on the ECD framework, beginning levels only involve area unit iteration skill and the other skills and related features are gradually added to the later game levels. Figure 3 shows the frequency of different events in each game level. As depicted in Figure 3, the student interactions with the system in all game levels is dominated by movements.

We only used the trajectories of the students who participated in the post-test. In the case of multiple attempts in a game

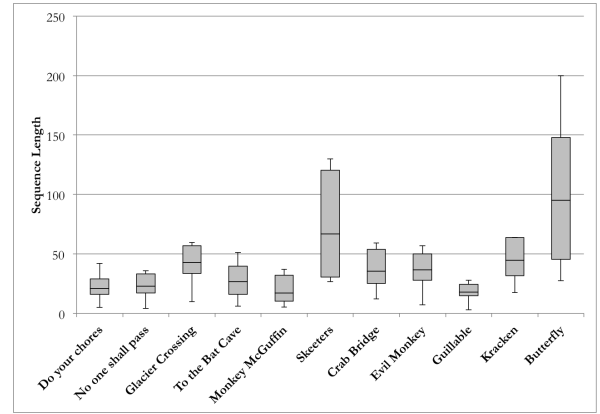


Figure 4. Boxplot of Sequence Length in Each Game Level

level, we only considered the trajectory from the first attempt. Figure 4 shows the boxplot of sequence length in each game level.

Experimental Setup

First we divide students into two groups, one (80%) for training and development purposes and the other (20%) for test and verification. In the synthesis phase, we transform log data of different nature (eg. movements, use of the tools, requests for hints) in each game level to observable values that can be used as evidence for learning. In the modeling phase, we use the observations to train two Hidden Markov Model (HMM)s that capture the sequence of actions for high- and low-performing students. In the aggregation phase, we use the likelihood of students' sequence of action in order to build a regression model that is predictive of the post-test results. Finally, we test the regression model on the held-out set (20%) and report the results.

Results

In order to evaluate our regression model we decided to compare its performance against two baselines: 1) A regression model that uses success and failure of students in each game level as feature **Write equation** and 2) A regression model that uses the sequence lengths in each game level **Write equation**. Our model was able to significantly outperform both baselines regarding to mean absolute error and root mean squared error. Moreover, the R^2 correlation and Spearman ρ of our predicted values with the true values was much higher than both baselines. Table 2 shows the results.

Predictive Features	R^2	ρ	MAE	RMSE
Sequence Length (Normal)	0.06	0.79	3.24	4.15
Success / Failure	0.01	0.63	3.22	4.14
SPRING	0.55	0.82	2.84	3.35

Table 2. The results of predicting post-test scores using three different feature sets

Show residual plot

Show confidence intervals

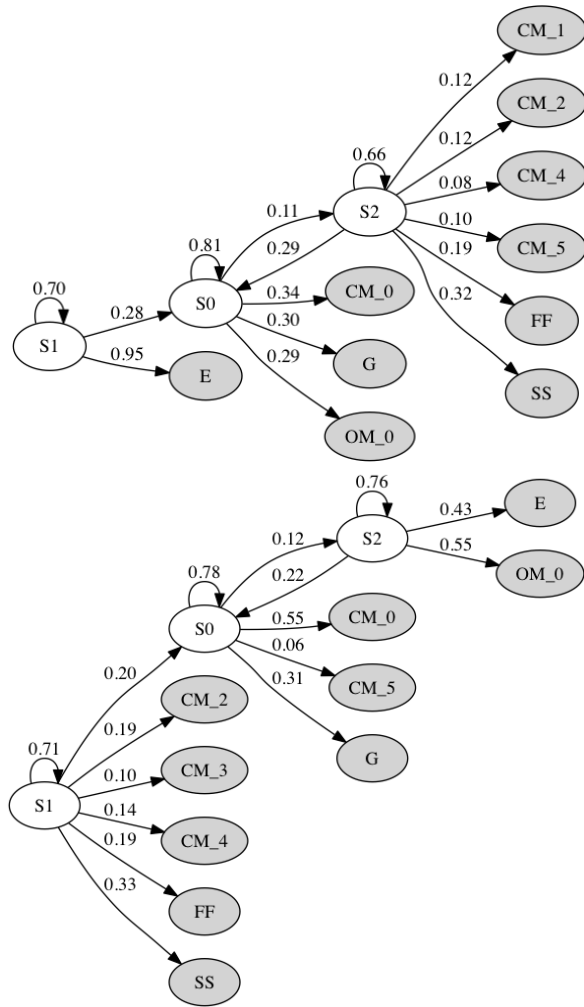


Figure 5. Two Hidden Markov Models learned from data in one of the game levels. The top HMM represents the high-performing students and the bottom one represents the low-performing ones.

Discussion

Figure 5 shows the difference between two HMMs learned for one of the game levels.

RELATION TO PRIOR WORK

In order to deal with such highly unstructured data, researchers often use carefully designed network structures (such as Bayesian Networks [1, 15]) or game-specific heuristics and benchmarks generated by experts playing the game [12, 17]. However, this approach is extremely labor intensive and might fail to capture meaningful patterns in student exploratory habits within the game. Given these limitations, data-driven analysis of student interactions provides a powerful alternative that facilitates the discussions around what does and does not work in a particular educational game.

The potential of computer games for educational purposes has been of interest since nearly the beginning of videogames. Unlike video games, which focus on creating an entertaining experience for the user, educational games require principles

and strategies that engage students while maximizing their learning gain. Therefore, data-driven analysis of student behavior is crucial to better understand the learning process and improve the tools in the future.

There have been numerous attempts among the educational research community to develop analytic methods and build predictive models based on the data from educational games. *Newton's Playground* is an ECD based educational game with 74 problems that designed to teach qualitative physics to students in eighth- and ninth-grade. Students have to guide a green ball to a red ball by creating simple machines. Everything obeys the basic rules of physics relating to gravity and Newton's three laws of motion. Shute et al. [15] studied the effect of ECD design on student learning and found that students who played the game in a 4 hour session, showed significant improved in their qualitative, conceptual physics understanding.

Rumble Blocks is another educational game designed to teach basic concepts of structural stability and balance to children in grades K-3 (ages 5-8 years old). Harpstead et al. [6], studied the alignment of game to its target learning goals by examining whether student solutions follows the targeted principals. They employ clustering techniques on the individual solutions created by actual students and use principle-relevant metric (PRM) to measure how closely the representative solution embodies a specific targeted principle. The results demonstrated a misalignment between the feedback provided to students and the targeted knowledge.

Battleship Numberline is another educational game for understanding fraction using number line estimation. Students attempt to explode target ships and submarines by estimating numbers on a number line. Lomas et al. [11] performed a large-scale online experiment in order to study the effect of challenge on player motivation and learning. They presented different configurations of the game for different groups of students and used a combination of time spent and challenges attempted as a measure of engagement and the average success rate of each design configuration as a measure of challenge. The results showed a linear correlation between challenge (difficulty) and engagement, meaning the easier the game, the longer students played.

Refraction is another educational game for learning about fractions by splitting laser beams into fractional amounts to target spaceships by avoiding asteroids. Liu et al. [10], created an ensemble algorithm that combines elements of Markov models, player heuristic search, and collaborative filtering techniques with state-space clustering in order to predict player movements on last game-level based on the history of movements in previous game levels. Lee et al. [9] extended the former framework by building state-action graph and using feature selection techniques to reduce the number of features for each state. To ensure extensibility, they also tested the framework on another game *DragonBox* and reported improvement over a Markov predictor.

CONCLUSION

Limitation, we don't compare to Valerie Schultz

Limitation, we don't control for game ability like Valerie does

Limitation, we only used data from Alice in wonderland

Method is accurate

REFERENCES

1. David W Albrecht, Ingrid Zukerman, and An E Nicholson. 1998. Bayesian models for keyhole plan recognition in an adventure game. *User modeling and user-adapted interaction* 8, 1-2 (1998), 5–47.
2. Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. 1996. A density-based algorithm for discovering clusters in large spatial databases with noise.. In *Kdd*, Vol. 96. 226–231.
3. Emily B Fox, Erik B Sudderth, Michael I Jordan, and Alan S Willsky. 2008. An HDP-HMM for systems with state persistence. In *Proceedings of the 25th international conference on Machine learning*. ACM, 312–319.
4. José P. González-Brenes and Jack Mostow. 2011. How to Classify Tutorial Dialogue? Comparing Feature Vectors vs. Sequences. In *Proceedings of the 4th International Conference on Educational Data Mining* (Eindhoven, Netherlands), M. Pechenizkiy and T. Calders and C. Conati and S. Ventura and C. Romero and J. Stamper (Ed.). 169–178.
http://educationaldatamining.org/EDM2011/wp-content/uploads/proc/edm2011_paper23_full_Gonzalez-Brenes.pdf
5. Mark Hall, Eibe Frank, Geoffrey Holmes, Bernhard Pfahringer, Peter Reutemann, and Ian H Witten. 2009. The WEKA data mining software: an update. *ACM SIGKDD explorations newsletter* 11, 1 (2009), 10–18.
6. Erik Harpstead, Christopher J MacLellan, Vincent Aleven, and Brad A Myers. 2014. Using extracted features to inform alignment-driven design ideas in an educational game. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, 3329–3338.
7. Peter Hofman, Bryan Goodwin, and Stuart Kahl. Re-Balancing Assessment: Placing Formative and Performance Assessment at the Heart of Learning and Accountability. (????).
8. M Lazarín. 2014. Testing overload in American schools. *Center for American Progress*. Retrieved from www.americanprogress.org/issues/education/report/2014/10/16/99073/testing-overload-in-american-schools (2014).
9. Seong Jae Lee, Yun-En Liu, and Zoran Popovic. 2014. Learning Individual Behavior in an Educational Game: A Data-Driven Approach. In *Educational Data Mining 2014*.
10. Yun-En Liu, Travis Mandel, Eric Butler, Erik Andersen, Eleanor O'Rourke, Emma Brunskill, and Zoran Popovic. 2013. Predicting player moves in an educational game: A hybrid approach. In *Educational Data Mining 2013*.
11. Derek Lomas, Kishan Patel, Jodi L Forlizzi, and Kenneth R Koedinger. 2013. Optimizing challenge in an educational game using large-scale design experiments. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, 89–98.
12. Travis Mandel, Yun-En Liu, Sergey Levine, Emma Brunskill, and Zoran Popovic. 2014. Offline policy evaluation across representations with applications to educational games. In *Proceedings of the 2014 international conference on Autonomous agents and multi-agent systems*. International Foundation for Autonomous Agents and Multiagent Systems, 1077–1084.
13. Robert J Mislevy, John T Behrens, Kristen E Dicerbo, and Roy Levy. 2012. Design and discovery in educational assessment: evidence-centered design, psychometrics, and educational data mining. *JEDM-Journal of Educational Data Mining* 4, 1 (2012), 11–48.
14. Valerie J Shute, Eric G Hansen, and Russell G Almond. 2008. You can't fatten a hog by weighing it-Or can you? Evaluating an assessment for learning system called ACED. *International Journal of Artificial Intelligence in Education* 18, 4 (2008), 289.
15. Valerie Jean Shute and Matthew Ventura. 2013. *Stealth assessment: Measuring and supporting learning in video games*. MIT Press.
16. Valerie J Shute, Matthew Ventura, Malcolm Bauer, and Diego Zapata-Rivera. 2009. Melding the power of serious games and embedded assessment to monitor and foster learning. *Serious games: Mechanisms and effects* 2 (2009), 295–321.
17. Bulent Tastan and Gita Reese Sukthankar. 2011. Learning Policies for First Person Shooter Games Using Inverse Reinforcement Learning.. In *AIIDE*.