

A Generalized Multidimensional Evaluation Framework for Player Goal Recognition

Wookhee Min,¹ Alok Baikadi,² Bradford Mott,¹ Jonathan Rowe,¹
Barry Liu,¹ Eun Young Ha,³ James Lester¹

¹Department of Computer Science, North Carolina State University, Raleigh, NC 27695

²Learning Research and Development Center, University of Pittsburgh, Pittsburgh, PA 15213

³IBM, Littleton, MA 01460

¹{wmin, bwmott, jprowe, bliu12, lester}@ncsu.edu, ²baikadi@pitt.edu, ³eyha@us.ibm.com

Abstract

Recent years have seen a growing interest in player modeling, which supports the creation of player-adaptive digital games. A central problem of player modeling is goal recognition, which aims to recognize players' intentions from observable gameplay behaviors. Player goal recognition offers the promise of enabling games to dynamically adjust challenge levels, perform procedural content generation, and create believable NPC interactions. A growing body of work is investigating a wide range of machine learning-based goal recognition models. In this paper, we introduce GOALIE, a multidimensional framework for evaluating player goal recognition models. The framework integrates multiple metrics for player goal recognition models, including two novel metrics, *n-early convergence rate* and *standardized convergence point*. We demonstrate the application of the GOALIE framework with the evaluation of several player goal recognition models, including Markov logic network-based, deep feedforward neural network-based, and long short-term memory network-based goal recognizers on two different educational games. The results suggest that GOALIE effectively captures goal recognition behaviors that are key to next-generation player modeling.

Introduction

While open-world digital games aim to promote players' engagement and replayability by enabling players to explore and achieve goals within expansive virtual worlds, the high degree of autonomy granted to players is often at odds with the task of game designers who have to provide coherent storylines and gameworld events (Riedl and Bulitko 2013; Min et al. 2014). To address the demands for player modeling, a broad range of computational approaches have been investigated for modeling players'

cognitive, affective and behavioral states in digital games (Yannakakis et al. 2013). Player modeling can support personalized game experiences in the context of interactive narrative (Riedl and Bulitko 2013), game balancing (Lopes and Bidarra 2011), and procedural content generation (Shaker et al. 2015).

Goal recognition, a central player modeling functionality, aims to dynamically identify high-level objectives that a player is attempting to achieve based on observable gameplay behaviors (Mott et al. 2006; Harrison et al. 2015). Open-world digital games pose a significant challenge for goal recognition (Min et al. 2016): these games do not explicitly provide a specific sequence of goals to complete in the game and there are a vast number of possible ways in which players can achieve a goal. Especially, in situations where players have limited prior experience with an open-world digital game, it may be the case that they perform exploratory actions instead of goal-directed actions, identify goals through a sequence of events that they have encountered and environments that they have observed, and then pursue a goal. This characteristic of open-world digital games results in highly idiosyncratic action sequences, and thus the task of recognizing players' goals exhibits significant uncertainty. Devising reliable computational models is key to the success of goal recognition in open-world digital games.

In this paper, we present a generalized multidimensional evaluation framework for goal recognition, GOALIE (Generalized Observable Action Learning for Intent Evaluation). The objective of this framework is to identify reliable goal recognition models that should not only correctly predict player goals on overall sequences of actions, but also make early predictions (i.e., making consistently correct predictions as early as possible) since run-time game adaptation is a central objective of goal recognition. GOALIE supports multidimensional evaluations of goal

recognition models through five metrics including three conventional metrics (accuracy rate, convergence point, and convergence rate) (Blaylock and Allen 2003) and two novel metrics (standardized convergence point and n -early convergence rate).

Our previous work (Min et al. 2014) found problems with the conventional convergence point metric: it can be misleading with respect to favoring goal recognition models with lower convergence rates. GOALIE addresses this problem by introducing two novel metrics. To illustrate goal recognition model evaluation, we use GOALIE to analyze three competitive computational goal-modeling techniques: long short-term memory networks (LSTMs) (Min et al. 2016), n -gram encoded feedforward neural networks pre-trained with stacked denoising autoencoders (Min et al. 2014), and Markov logic networks (Baikadi et al. 2014) on two goal recognition data corpora from two independent open-world educational games with data from 137 students and 828 students, respectively.

Results suggest that GOALIE can be used to identify the goal recognition models that are most promising for player modeling. For example, GOALIE finds that LSTM-based goal recognition harnessing distributed action embeddings significantly outperforms the other two approaches, with respect to predictive accuracy, convergence rate, standardized convergence point, and n -early convergence rate, across the two benchmark datasets. Even though the LSTM-based approach does not achieve the highest performance in terms of the conventional convergence point, its complementary metric, standardized convergence point, suggests that the LSTM-based approach outperforms the competitive baseline approaches with respect to the models' early prediction capacity.

Related Work

Because plan recognition focuses on inferring plans and goals of observed agents, it is formulated as a generalized task of goal recognition (Baker et al. 2009; Sukthankar et al. 2014). While much previous plan recognition work has utilized a hand-crafted plan library and a decision model (Geib and Goldman 2009; Bisson et al. 2015), a salient line of investigation has addressed plan recognition by learning a plan library in a data-driven approach (Fagan and Cunningham 2003; Synnaeve and Bessière 2011) or dispensing with the need for a plan library (Baker et al. 2009; Ramírez and Geffner 2011) by interpreting plan recognition as an inversion of action planning given a goal. However, in open-world digital games, particularly those in which players have little or no prior experience, players' exploration-based actions are marked by highly idiosyncratic sequences of player actions and often sub-optimal for achieving goals (Min et al. 2016); thus, devising a reli-

able plan library or using a planning approach that assumes a rational agent is infeasible. Corpus-based, statistical goal recognition (Blaylock and Allen 2003) effectively deals with this challenge in open-world digital games by requiring only a list of goals and a corpus containing action sequences that achieve the goals, and thus holds significant potential for performing high-level game adaptations.

Along with deep learning's significant advance in computer vision, speech recognition and natural language processing (LeCun et al. 2015), it has demonstrated considerable success in goal and plan recognition (Bisson et al. 2015; Min et al. 2016), perhaps because of the focus on extracting hierarchical representations from lower-level inputs (i.e., actions) to higher-level outputs (i.e., goals) (Min et al. 2014). Bisson and colleagues (2015) investigated recursive neural network-based decision models and evaluated the approach on plan recognition domains including a real-time strategy game. Min et al. (2016) examined goal recognition with LSTMs, which achieved state-of-the-art predictive accuracy by effectively modeling player behavior sequences in an open-world game. In this paper, we illustrate the application of GOALIE with the LSTM-based goal modeling technique along with two competitive baseline approaches on two different data corpora.

GOALIE Framework

A standard framework and set of metrics for goal recognition have not been established (Sukthankar et al. 2014). A commonly used approach measures models' predictive performance such as their accuracy rate (Ha et al. 2011). As a supplementary set of metrics, Blaylock and Allen (2003) proposed the convergence point and convergence rate that capture models' early prediction capacity, which has been investigated in a wide range of goal recognition work (Mott et al. 2006; Baikadi et al. 2014). Early prediction is of significant importance because goal recognizers that lack an early prediction capacity cannot effectively support goal-driven gameplay adaptation at run-time.

In this work, we present GOALIE, a generalized multi-dimensional evaluation framework, equipped with two novel complementary metrics: *standardized convergence point* and *n -early convergence rate* as well as the conventional metrics. These two new metrics are associated with the convergence point and rate, and suggest reinterpreted results with respect to goal recognition models' early prediction capacity. In the subsequent paragraphs, we describe how the new metrics featured in GOALIE differ from the conventional convergence metrics.

The metric of convergence point (Blaylock and Allen 2003) measures how early goal recognition models can consistently make accurate predictions within a converged

sequence, i.e., an action sequence in which the last goal prediction is correct. More formally, convergence point is calculated by $\sum_{i=1}^m (k_i/n_i) / m$, in which m is the number of converged action sequences, and n_i and k_i are the total number of actions and the number of actions after which the goal recognizer consistently makes accurate predictions in the i^{th} converged action sequence, respectively (Min et al. 2016). Note that convergence point ignores all action sequences that do not converge to a correct goal.

In this paper, we revisit the convergence point metric, which provides misleading results with respect to early prediction in some cases. For example, suppose that we have two goal recognition models: GR₁ and GR₂, and there are two action sequences (AS₁ and AS₂) to predict:

- AS₁: Action₁₁, Action₁₂, and Action₁₃, whose goal is G_A.
- AS₂: Action₂₁, Action₂₂, and Action₂₃, whose goal is G_B.

In this situation, assume that both GR₁ and GR₂ correctly predict the goal (G_A) associated with the three actions in AS₁, while GR₁ makes incorrect predictions on all three actions of AS₂, but GR₂ makes correct predictions for the goal of the last two actions (Action₂₂ and Action₂₃). The convergence point of GR₁ is 0.33 since it consistently makes correct predictions after observing the first action for AS₁, which is the only converged sequence. However, the convergence point of GR₂ is 0.5 which is computed as $(1/3+2/3)/2$, and thus GR₁ is identified as the model with better early prediction capacity based on convergence point (lower is better and the maximum possible value is one for this metric) even though GR₂ accurately predicts actions sooner than GR₁ in AS₂. This issue is exacerbated for a majority class-based goal recognition model that predicts a single goal for all actions. This model, which is unreliable, will yield a low convergence point by making correct predictions for all actions in converged action sequences. To overcome this issue with the conventional convergence point, we introduce *standardized convergence point*.

Standardized convergence point

The standardized convergence point metric measures a convergence point regardless of whether an action sequence converged to a correct goal or not. To compute this metric, a non-converged action sequence has a convergence point of (the total number of actions + penalty) divided by (the total number of actions), thereby yielding a convergence point greater than one. In this manner, non-converged action sequences are penalized in terms of early prediction.

Definition 1. Standardized convergence point is calculated by $\sum_{i=1}^m (k_i/n_i) / m$, in which m is the total number of action sequences, and n_i is the total number of actions in the i^{th} action sequence. k_i is contingent on whether the i^{th} action sequence converged or not; if converged, k_i is the number of actions after which the goal recognizer consistently

makes accurate predictions as in the conventional convergence point metric; otherwise, k_i is n_i+p_i , where p_i ($p_i > 0$) is the penalty parameter for the i^{th} action sequence. Lower is better for this metric.

Returning to the example presented in the preceding subsection, standardized convergence point (in this work, we set p_i to 1 for all action sequences) of GR₁ is 0.83, computed by $(1/3+4/3)/2$, while standardized convergence point of GR₂ remains the same as 0.5. Thus, using this definition, the standardized convergence point of GR₂ is lower than GR₁. Since early prediction is more accurately captured using the standardized convergence point than the convergence point as discussed, GOALIE measures models' early prediction capacity with this new metric.

N-Early Convergence Rate

Convergence rate (Blaylock and Allen 2003) measures the percentage of action sequences in which the last goal prediction is correct. While the goal recognition system generally makes correct predictions on individual actions (i.e., high accuracy rate), if it makes an incorrect prediction on the last action (i.e., low convergence rate), the model lacks reliability since incorrect game adaptation at the end of an action sequence might confuse players who are about to achieve the planned goal. We extend this metric to consider the last $(n+1)$ action predictions through a new metric, *n-early convergence rate*.

Definition 2. *N*-early convergence rate is calculated by $\sum_{i=1}^m k_i / m$, in which m is the total number of action sequences. k_i is 1 if the last $(n+1)$ goal predictions are all correct in the i^{th} action sequence; otherwise, k_i is 0 for the action sequence. In a special case of this metric when n equals 0, the definition of *n*-early convergence rate is the same as the conventional convergence rate. Higher is better for this metric, in which the maximum value is one.

N-early convergence rate views goal recognizers' early prediction from a different perspective compared to standardized convergence point. The *n*-early convergence rate takes a static approach (i.e., a fixed window size of $n+1$) backward from the end of an action sequence, while standardized convergence point takes a dynamic approach forward from the beginning of an action sequence. These two novel measurements are designed to complement the corresponding conventional convergence metrics when evaluating goal recognition models' early prediction capacity.

It should be noted that neither of these novel convergence metrics nor the conventional convergence metrics can be dynamically computed during gameplay since players' goals are hidden from the goal recognizer. Nonetheless, offline measurements of these metrics offer considerable insight into selecting the most reliable goal recognition model to best support runtime game adaptation.

Goal Recognition Corpora

In this section, we describe two different educational open-world digital games featuring middle-grade science education (CRYSTAL ISLAND: OUTBREAK) and elementary-grade science education (CRYSTAL ISLAND: UNCHARTED DISCOVERY). Goal recognition corpora were collected from players' trace logs playing the two educational games, which allow us to perform post-hoc data analyses of goal recognition.

As noted above, players' action sequences do not necessarily represent optimal paths for achieving goals in both games since the players did not have prior experience playing them. Rather, players often explore the virtual environment in order to familiarize themselves with the game-world, and make gradual but possibly circuitous progress toward each objective, eventually culminating in the final action that achieves a goal (Min et al. 2016).

CRYSTAL ISLAND: OUTBREAK (Rowe et al. 2011) is an open-world educational game for middle school science. In the game players are tasked with identifying the cause of an illness afflicting a team of scientists on a remote island research camp. Implemented in the Half-Life 2 Source engine, its gameplay is similar to many exploration-centric games in which players perform actions such as navigating, discovering important items, and talking with non-player characters. The non-linear narrative consists of seven key goals that players must accomplish to complete the game. Five of these goals involve speaking with NPCs about the spreading illness, while the remaining two involve testing contaminated food in the camp's laboratory and submitting a completed diagnosis to the camp nurse.

The data were collected during a 60-minute-long study involving 137 eighth grade students from a public middle school. The CRYSTAL ISLAND: OUTBREAK data corpus consists of 77,182 player actions (i.e., the total number of possible goal recognitions) and 893 achieved goals, with an average of 86.4 player actions per goal.

CRYSTAL ISLAND: UNCHARTED DISCOVERY (Baikadi et al. 2014; Lester et al. 2014) is an open-world educational game for upper elementary science. The game takes place on a fictional island in the Oceania region of the Pacific Ocean. The player joins a cast of virtual characters on the island in establishing a new life, and explores an expansive virtual world as she performs tasks for the island inhabitants. As the player learns about landforms, navigation and modeling, she is asked to perform several quests that assess their problem-solving skills.

In the dataset used in this work, players completed four quests. Two of the quests focused on understanding landforms, such as plateaus, deltas and waterfalls. The other two quests involved understanding navigation, both through reading a map and following a heading for a specified distance. Each quest was associated with three goals.

Once a quest has begun, players can pursue its three goals in any order, and also can initiate multiple quests at the same time. The total number of possible goals is 12.

The data used in the evaluation of our goal recognition model was collected during a four-week study involving 828 fifth grade students from eight public elementary schools. The dataset consists of 811,647 player actions and 7,652 achieved goals, with an average of 106.1 player actions per goal.

Long Short-Term Memory Network-Based Goal Recognition Model

We cast goal recognition as a multiclass classification problem in which a trained classifier predicts the most likely goal associated with the currently observed action sequence (Min et al. 2016). We assume that a given sequence of actions maps to a single goal, and no interleaving occurs between actions associated with different goals, since the goal recognition corpus does not lend itself to this type of analysis (Ha et al. 2011).

To devise a reliable goal recognition model that effectively deals with noisy player actions, our preliminary work has focused on the exploratory nature in open-world digital games. Player goals and actions constitute cyclical relationships (Ha et al. 2011); players' previously achieved goals may inform their subsequent actions, and their current actions may influence their upcoming goals. Consequently, modeling sequential patterns from a series of player actions and previously achieved goals is crucial for accurate goal recognition (Min et al. 2016). Inspired by these characteristics, our previous work investigated long short-term memory networks (LSTMs) (Min et al. 2016), n -gram encoded feedforward neural networks pre-trained with stacked denoising autoencoders (Min et al. 2014), and Markov logic networks (Ha et al. 2011; Baikadi et al. 2014). In the remainder of this section, we present a brief overview of LSTM-based goal recognition in UNCHARTED DISCOVERY. Details regarding LSTM-based goal recognition in OUTBREAK can be found in (Min et al. 2016).

LSTMs are a type of recurrent neural networks (RNNs) that are specifically designed for sequence labeling of temporal data (Graves 2012). LSTMs have achieved high predictive performance in various sequence labeling tasks, often outperforming standard RNNs by preserving longer-term dependencies than standard RNNs and effectively addressing the vanishing gradient problem that occurs when training standard RNNs.

LSTMs (Figure 1) feature a sequence of memory blocks that include one or more self-connected memory cells along with three gating units: an input gate, a forget gate, and an output gate. In LSTMs, the input and output gates modulate the incoming and outgoing signals to the memory

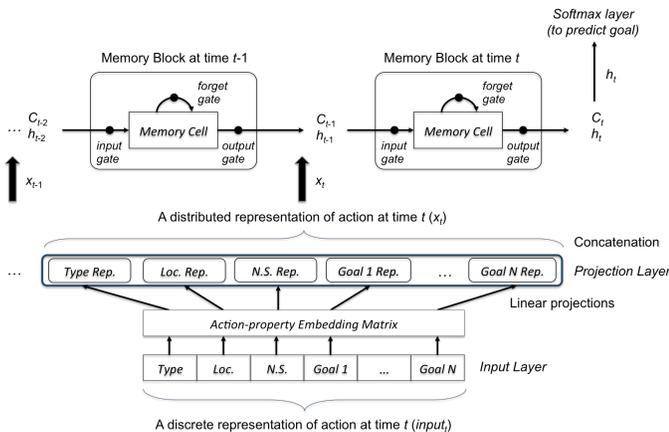


Figure 1. LSTM-based goal recognition (Min et al. 2016).

cell, and the forget gate controls whether the previous state of the memory cell is remembered or forgotten.

A player action is encoded with four properties: *action type*, *action location*, *narrative state* (the player’s progress in solving the game’s narrative), and *previously achieved goals* (Min et al. 2014). For UNCHARTED DISCOVERY, action types include 126 distinct types of player actions under 19 high-level action categories, action location includes 84 non-overlapping sub-locations within the gameworld, narrative states contain 16 possible values based on the interactive storyline’s plot structure, and there are 13 possible goals that could be previously achieved, including ‘None’ in case the player has not yet achieved any goals.

As depicted in Figure 1, the input layer is fed with a $(N+3)$ -dimensional discrete vector (we set N to 12 in this work). The first three dimensions of the vector are allocated to represent the action type, action location, and current narrative state with integer-based indices, while the following N dimensions represent a sequence of previously achieved goals also with integer-based indices.

The next layer is a projection layer that creates dense distributed vector representations out of discrete representations in the input layer (Bengio et al. 2003; Min et al. 2016). To support the projection layer, a shared action-property embedding matrix, the size of which is 239 (the sum of possible values of action properties computed as $126+84+16+13$) by d (embedding size), is created (Figure 1). The action-property embedding matrix is randomly initialized following a uniform distribution and is fine-tuned during supervised machine learning.

At prediction time, once a distributed representation per action property of an action at time t is retrieved from the embedding matrix, these vector representations are concatenated into a single $(N+3)*d$ dimensional vector, which is fed into the LSTM model as an input (x_t) at time t . The final memory cell output vector (h_t in Figure 1) is used to

predict the most likely goal for the sequence of actions in a softmax layer.

For LSTM model configurations, we explore one hyperparameter: the number of hidden units between 100 and 200, while fixing other hyperparameters. This hyperparameter was found to make the most significant influence on the goal recognition model performance for OUTBREAK (Min et al. 2016), and we investigate it with higher priority over other hyperparameters. We use the embedding size (d) of 20 and the dropout rate of 0.75 (Srivastava et al. 2014), adopt a mini-batch gradient descent with the mini-batch size of 128, and utilize categorical cross entropy for the loss function and the Adam stochastic optimizer (Kingma and Ba 2015). For training efficiency, action sequences longer than ten are pruned to keep only the last ten actions. Finally, the training process stops early if the validation score has not improved within the last seven epochs. In each fold, 10% of the training data is used to determine early stopping, and 90% is utilized for supervised training, while the validation data is purely used for calculating the validation score. The maximum number of epochs is set to 100.

Evaluation

We illustrate the GOALIE framework by applying it to multiple goal recognition models for the two CRYSTAL ISLAND game-based learning environments. For OUTBREAK, the distribution of the 7 goals ranges from 6.4% (“Speaking with the camp nurse”) to 26.6% (“Running laboratory test on contaminated food”), which is the majority class baseline. The UNCHARTED DISCOVERY data corpus consists of 12 goals, in which the most frequent goal (“Picking up the dark blue flag”) and the least frequent goal (“Placing a sign at volcano”) appeared 13.5% (i.e., majority class baseline) and 3.1% out of the entire set of players’ achieved goals, respectively.

For n -gram encoded feedforward neural networks (FFNNs), we explore two hyperparameters: the number of hidden units among $\{100, 200\}$ and the number of hidden layers among $\{2, 3\}$, while using a fixed value for other hyperparameters as follows: n of 5 for n -gram encoding and the corruption level of 0.5 for stacked denoising auto-encoders (these values were found to be the best configurations for OUTBREAK (Min et al. 2014)). For FFNNs, we adopt stochastic gradient descent for both unsupervised pre-training and supervised fine-tuning, with learning rates of 0.1 for pre-training and 1 for fine-tuning.

Through GOALIE, we compare three goal-modeling techniques, including LSTMs, FFNNs, and Markov logic networks (MLNs) (Domingos et al. 2006) on the two data corpora, and identify the most reliable goal recognition model for each data corpus. We evaluate induced models

using three conventional metrics: accuracy rate, convergence rate, and convergence point as well as the two proposed novel metrics, standardized convergence point and n -early convergence rate. Due to space limitations, we only report 1-early convergence rate for n -early convergence rate, by which we evaluate whether goal recognizers correctly predict the goal for the last two actions in every action sequence. While LSTMs and FFNNs have identified the best performing models through a machine-driven approach (i.e., grid search of hyperparameters based on cross validation results), MLNs have utilized human expert-crafted logic formulae in terms of discovery events, domain-specific representations of user progress specifically targeted to each digital game (Baikadi et al. 2014). We use the same data split in 10-fold cross validation for fair comparisons.

Table 1. Averaged rates of MLNs, FFNNs, and LSTMs for CRYSTAL ISLAND: OUTBREAK.

	MLN	FFNN	LSTM
Accuracy Rate (%)	55.21	62.43	66.35
Convergence Point (%)	30.80	41.30	33.34
Stand. Convergence point (%)	67.66	62.66	53.19
Convergence Rate (%)	49.09	70.06	71.32
1-Early Convergence Rate (%)	46.71	64.93	68.81

Table 1 presents results of the three computational goal recognition approaches for the OUTBREAK data corpus. In this table, only the model that achieves the highest cross-validation accuracy per approach is reported, which are 100-hidden-unit models for LSTMs and 2-hidden-layer models with 100 hidden units per layer for FFNNs. LSTMs achieve the best goal recognition accuracy (66.35%), convergence rate (71.32%), 1-early convergence rate (68.81%), and standardized convergence point (53.19%), while MLNs achieve the best (i.e., lowest) convergence point (30.80%). As noted above, the convergence point is calculated only for converged sequences (i.e., 49.09% of the total action sequences for MLNs), and the result indicates that MLNs achieve the most efficient early prediction for a relatively smaller number of action sequences. LSTMs and FFNNs overall produce more reliable predictions on more action sequences than MLNs as evidenced by higher accuracy rates and convergence rates, but the high convergence rates led the models to consider more noisy action sequences that eventually converged to correct goals but are not trivial to predict, thereby inducing higher convergence points. This phenomenon has been called *inherent tension between convergence rate and convergence point* in previous work (Min et al. 2014). The standardized convergence point suggests reinterpreted results on the early prediction. LSTMs achieve the lowest standardized convergence point with sizable difference over the other two approaches. Based on these results, we conclude that

LSTM is the best goal modeling technique for the CRYSTAL ISLAND: OUTBREAK game.

Table 2 presents results for UNCHARTED DISCOVERY following the same process described in OUTBREAK. For this data corpus, the best model configurations for LSTMs and FFNNs are 200-hidden-unit models and 2-hidden-layer models with 200 hidden units per layer, respectively.

Similar to the results for OUTBREAK, GOALIE suggests that LSTM-based goal recognition models achieve the best performance for UNCHARTED DISCOVERY in every evaluation metric except for convergence point. The standardized convergence point modulates the inherent tension between convergence rate and convergence point for the models, and therefore GOALIE concludes that LSTM is the best goal modeling technique for this dataset as well.

Table 2. Averaged rates of MLNs, FFNNs, and LSTMs for CRYSTAL ISLAND: UNCHARTED DISCOVERY.

	MLN	FFNN	LSTM
Accuracy Rate (%)	24.40	32.26	35.32
Convergence Point (%)	75.54	47.38	53.66
Stand. Convergence Point (%)	91.71	75.73	75.63
Convergence Rate (%)	43.64	49.94	56.15
1-Early Convergence Rate (%)	39.05	44.56	49.19

Conclusion

Goal recognition is a core player modeling functionality in open-world digital games. We have introduced the GOALIE framework that performs multidimensional evaluation of goal recognition models. Empirical evaluations with GOALIE indicate that LSTMs achieve the most reliable results across all metrics on the two examined digital games, except for the convergence point. However, standardized convergence point, a novel early prediction metric proposed in GOALIE, suggests that LSTMs exhibit an improved early prediction capacity over FFNNs and MLNs, and thus, with GOALIE, we conclude that LSTM is the best performing goal modeling technique for the two open-world digital games.

In the future it will be important to investigate how to set the penalty parameter for the standardized convergence point and n in n -early convergence rate for goal recognition models targeting open-world digital games, and prioritize goal recognizers when having conflicting evaluation results with GOALIE. Moreover, an additional set of metrics that quantify unmeasured aspects of goal recognizer performance and a visualization tool that illustrates dynamic performance changes of goal recognition models across time would complement the current implementation of GOALIE. Finally, it will be important to investigate how goal recognition models identified through GOALIE support game adaptation at run-time.

References

- Baikadi, A., Rowe, J., Mott, B. and Lester, J. 2014. Generalizability of Goal Recognition Models in Narrative-Centered Learning Environments. In *Proceedings of the 22nd User Modeling, Adaptation, and Personalization*, 278–289. Springer International Publishing.
- Baker, C., Saxe, R. and Tenenbaum, J. 2009. Action understanding as inverse planning. *Cognition* 113(3): 329–349.
- Bengio, Y., Ducharme, R., Vincent, P. and Janvin, C. 2003. A Neural Probabilistic Language Model. *The Journal of Machine Learning Research* 3: 1137–1155.
- Bisson, F., Larochelle, H. and Kabanza, F. 2015. Using a Recursive Neural Network to Learn an Agent’s Decision Model for Plan Recognition. In *Proceedings of the 24th International Conference on Artificial Intelligence*, 918–924.
- Blaylock, N. and Allen, J. 2003. Corpus-based, statistical goal recognition. In *Proceedings of the 18th International Joint Conference on Artificial Intelligence*, 1303–1308.
- Domingos, P., Kok, S., Poon, H., Richardson, M. and Singla, P. 2006. Unifying Logical and Statistical AI. In *Proceedings of the 21st National Conference on Artificial Intelligence*, 2–7.
- Fagan, M. and Cunningham, P. 2003. Case-based plan recognition in computer games. In *Proceedings of the 5th International Conference on Case-Based Reasoning*, 161–170. Springer Berlin Heidelberg.
- Geib, C. and Goldman, R. 2009. A probabilistic plan recognition algorithm based on plan tree grammars. *Artificial Intelligence* 173(11): 1101–1132.
- Graves, A. 2012. *Supervised sequence labelling with recurrent neural networks*. Vol. 385. Springer.
- Ha, E. Y., Rowe, J. P., Mott, B. W. and Lester, J. C. 2011. Goal Recognition with Markov Logic Networks for Player-Adaptive Games. In *Proceedings of the 7th AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, 32–39.
- Harrison, B., Ware, S., Fendt, M. and Roberts, D. 2015. A Survey and Analysis of Techniques for Player Behavior Prediction in Massively Multiplayer Online Role-Playing Games. *IEEE Transactions on Emerging Topics in Computing* 3(2): 260–274.
- Kingma, D. P. and Ba, J. L. 2015. Adam: a Method for Stochastic Optimization. In *Proceedings of the 3rd International Conference on Learning Representations*.
- LeCun, Y., Bengio, Y. and Hinton, G. 2015. Deep Learning. *Nature* 521(7553): 436–444.
- Lester, J. C., Spires, H. A., Nietfeld, J. L., Minogue, J., Mott, B. W. and Lobene, E. V. 2014. Designing game-based learning environments for elementary science education: A narrative-centered learning perspective. *Information Sciences* 264: 4–18.
- Lopes, R. and Bidarra, R. 2011. Adaptivity challenges in games and simulations: A survey. *IEEE Transactions on Computational Intelligence and AI in Games* 3(2): 85–99.
- Min, W., Ha, E. Y., Rowe, J., Mott, B. and Lester, J. 2014. Deep Learning-Based Goal Recognition in Open-Ended Digital Games. In *Proceedings of the 10th AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, 37–43.
- Min, W., Mott, B., Rowe, J., Liu, B. and Lester, J. 2016. Player Goal Recognition in Open-World Digital Games with Long Short-Term Memory Networks. In *Proceedings of the 25th International Joint Conference on Artificial Intelligence*, 2590–2596.
- Mott, B., Lee, S. and Lester, J. 2006. Probabilistic goal recognition in interactive narrative environments. In *Proceedings of the 21st National Conference on Artificial Intelligence*, 187–192. Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999.
- Ramírez, M. and Geffner, H. 2011. Goal recognition over POMDPs: Inferring the intention of a POMDP agent. In *Proceedings of the 22nd International Joint Conference on Artificial Intelligence*, 2009–2014.
- Riedl, M. O. and Bulitko, V. 2013. Interactive Narrative: An Intelligent Systems Approach. *AI Magazine* 34(1): 67–77.
- Rowe, J., Shores, L., Mott, B. and Lester, J. 2011. Integrating learning, problem solving, and engagement in narrative-centered learning environments. *International Journal of Artificial Intelligence in Education* 21(1–2): 115–133.
- Shaker, N., Togelius, J. and Nelson, M. J. 2015. *Procedural Content Generation in Games: A Textbook and an Overview of Current Research*. Springer.
- Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I. and Salakhutdinov, R. 2014. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research (JMLR)* 15:1929–1958.
- Sukthakar, G., Geib, C., Bui, H., Pynadath, D. and Goldman, R. P. 2014. *Plan, Activity, and Intent Recognition*. Elsevier.
- Synnaeve, G. and Bessière, P. 2011. A Bayesian model for opening prediction in RTS games with application to StarCraft. In *2011 IEEE Conference on Computational Intelligence and Games*, 281–288. IEEE.
- Yannakakis, G., Spronck, P., Loiacono, D. and Andre, E. 2013. *Player Modeling*. In Lucas, S. M., Mateas, M., Preuss, M., Spronck, P. and Togelius, J. (eds.) *Dagstuhl Follow-Ups*. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, Dagstuhl, Germany.