

Methods and Standards for Research on Explainable Artificial Intelligence: Lessons from Intelligent Tutoring Systems

Robert Hoffman¹ and William Clancey¹

¹Florida Institute for Human and Machine Cognition

June 8, 2021

Abstract

We reflect on the progress in the area of Explainable AI (XAI) Program relative to previous work in the area of intelligent tutoring systems (ITS). A great deal was learned about explanation—and many challenges uncovered—in research that is directly relevant to XAI. We suggest opportunities for future XAI research deriving from ITS methods, as well as the challenges shared by both ITS and XAI in using AI to assist people in solving difficult problems effectively and efficiently.

INTRODUCTION

The general objective of XAI is to develop methods that enable practical use of an AI tool, including understanding the system’s capabilities and vulnerabilities. This knowledge makes it possible for users to act appropriately, such as cross-checking and complementing the automated work to accomplish the intended function within a broader established activity. “Explanation” is one way to assist people in gaining this expertise.

What are the best ways to explain complex systems? Can we facilitate learning by promoting self-explaining? What pedagogical approaches should computer-based tutoring systems use, and should they be derived from studies teachers interacting with students? These were among the questions driving AI research in the area of Intelligent Tutoring Systems (ITSs) since the 1970s [23, 24, 27]. We illustrate this work with an ITS for image interpretation that uses statistical analysis to relate features of images, MR Tutor [22]. In MR Tutor “explanation” is framed as an *instructional activity* for learning how to carry out a diagnostic task using an AI program as an aid. The following sections outline how models in this program and other ITS systems are created and used, followed by comparison to XAI objectives and methods.

OVERVIEW OF INTELLIGENT TUTORING SYSTEMS

Intelligent Tutoring Systems research—more broadly known as the field of “AI and Education” [1]—has been concerned with individualized instruction in many different domains, with a variety of representational media and interactive methods [27]. It should be mentioned at the outset that ITS systems work, and they can work quite well at teaching STEM topics compared to teacher-to-student tutoring [1]. Furthermore, it is important to note that ITS systems generally have not been intended as replacements for human teachers or tutors, but rather designed as tools to assist in classwork and for independent learning.

ITS work can be traced to the 1960s with the development of AI programs that represent knowledge in *structured models*, especially semantic nets, production rules, and schemas/frames. Programs using such models can solve problems, such as answering factual questions, proving theorems, and manipulating mathematical equations. Subsequent research in the 1970s added a reasoning module that interprets the structured model in particular situations for professional tasks, such as diagnosis, planning, design, and process control; these programs were called “expert systems.”

In general, an *intelligent tutoring program* contains such an AI problem-solving program, using it to interact with and instruct a student. Thus ITS is contrasted with computer-based instruction programs that do not have a built-in capability to solve the problems that are presented to students.

Most intelligent tutoring programs engage the student in a learning activity in which the program serves as an instructor; they use distinct models of the domain, the student, and curriculum; and the interactive design is based on a theory of the pedagogical process [13, 23, 24, 27].

ITS research has been concerned with teaching mathematics [1] and basic science, as well as professional expertise relating to complicated systems, such as electro-mechanical troubleshooting [3], engineering operations [16], and medicine [9]. Insofar as machine learning programs are complicated systems whose capabilities and, to some extent, methods we want users to understand, the techniques and lessons from ITS research and development over nearly 50 years are worth considering for adoption in XAI research.

COMPARISON OF ITS AND XAI APPROACHES TO EXPLANATION

The most obvious distinction between ITS and XAI programs in general is that ITS programs manipulate and interpret three essential models: 1) the formalization of what is to be learned (*domain model*); 2) hypotheses about the user’s knowledge and cognitive activity (*student model*); and 3) a pedagogical method (e.g., applying a *curriculum model*).

Strikingly, XAI programs attempt to help users learn—which is after all the intent of providing explanations—without similar models and interactive capabilities. In effect, ITS and XAI researchers have proceeded with different perspectives on the learning activity, the nature of explanation, importance of understanding the user/student’s mental model and reasoning method, and following a comprehensive pedagogical plan for instructing the student. We consider each of these aspects briefly here.

Learning Activity and Instructional Setting

In general, a student is presented by an ITS program with some question or problem that he or she must solve—a form of “learning by doing.” A great deal of ITS research since 1990 has explored different kinds and aspects of a learning activity: the *setting* (e.g., “job-embedded” tutors; museum displays), *media* (e.g., web-based presentations; virtual worlds with role-playing simulations; animated agents), *instructional mode* (e.g., “collaborative inquiry”), *human-machine interaction* (e.g., recognizing and conveying emotion; speech understanding), *theories of learning*, a central theme in the “learning sciences” (e.g., cognitive processes, situated learning, apprenticeship), *educational data-mining* (e.g., learning patterns from web-based databases of student performance), and *supporting teachers* (e.g., calling attention to a student requiring remedial work).

Woolf [27] surveyed recent perspectives about ITS applications, emphasizing interactive approaches in which the student is not just a passive recipient of instruction, aka “active learning.” Clancey and Soloway [13] edited a journal, *Interactive Learning Environments*, which also emphasized theory-based design of an instructional setting and the role of viewing and manipulating media in the learning process.

The Nature of Explanation

XAI research generally assumed at first that “explanation” only involves the process of providing an explanation to the user, on the assumption that an explanation consists of text or a graphic, which is good and sufficient in itself. But ITS research clearly demonstrated how explanation must be understood from the user’s perspective as *a learning process*, and thus from the program’s perspective as an instructive process (which includes explaining) rather than a “one-off,” stand-alone question-answer interaction. This is true whether the learning process is an activity involving a person and machine, a group of people, or process of self-explanation by a person or program. For some XAI applications, explanation will be part of an activity that extends over multiple uses and interactions, especially because a neural network program can continually evolve. XAI researchers have thus begun to consider ways in which the user can explore how the AI program works and its vulnerabilities (a concern ignored by that ITS programs that focus on textbook knowledge).

ITS research demonstrated that “explanation” is an interaction among the user, the artifact, and their activity in a task context. In particular, the format/medium, content, and timing of explanations may differ to support different information needs for different tasks. In critical, time-pressed situations the only practical support may be directing the user’s attention; in activities over hours or days, such long-term care for a patient, the program may serve more as an assistant in constructing situation-specific models and action plans.

The process of instruction, including explaining, necessarily involves shared languages and methods for communicating. The earliest ITSs demonstrated some form of natural language capability, such as mixed-initiative question-answering, case-method dialogue, Socratic discourse, or customized narrative presentations. ITSs have also used graphic presentations and animated simulations to convey relationships and causality. Similarly, a general consensus has emerged among XAI researchers that the explanation process must involve the exchange of meaningful (versus computationally formal) information.

Domain Model

The domain model is a representation of the subject material to be taught. Generally, this model represents some process or system in a computationally interpretable form. This is a defining characteristic of “symbolic AI”; such models are variously described as qualitative, relational, and/or semantic. Clancey [7] presented evidence that all expert systems were “model-based.” In effect, a fundamental contribution of expert systems research has been to extend modeling formalisms beyond mathematical constructs in conventional science and engineering to include qualitative/relational models and associated modeling operations [10].

By design, an ITS domain model can be interpreted to solve problems presented to the student. This model constitutes the knowledge to be learned: facts, formalisms, causal processes, and reasoning processes (e.g., a diagnostic strategy). Typically, the program that applies the domain model to solve problems is called the *inference engine* (e.g., a rule interpreter), which applied to particular circumstances (e.g., a patient’s history and symptoms), produces a (partial) solution called a *situation-specific model*.

At first, most ITS researchers viewed the domain model as being isomorphic to stored structures and subconscious processes in the brain. This assumption motivated much productive research in which the ITS creates a model of the student’s knowledge and reasoning. In recent decades, models have been viewed more often as *scientific tools* that are constructed and applied to understand and manipulate processes and systems in the world [11, 21], exemplified by the student model in an ITS.

User (Student) Models

The student model represents the student’s knowledge relative to the domain model, in both general and situation-specific forms [8]. The student model may be constructed using the overlay method (student model is a subset of the expert domain model), the misconception/bug method (student behavior is matched against variants/incorrect domain model), or a machine learning method [17].

Some systems can also infer the student’s ongoing approach to solving a problem, such as a diagnostic strategy. ITS programs infer the student’s model by asking (e.g., “Do you believe X causes Y?”) or by interpreting reasoning steps. Misconceptions can also be pre-enumerated in the program and available for matching against student behavior [23]. The creation of models of student knowledge and reasoning remains an ongoing concern in ITS research.

Given the focus and challenge of adding an explanation capability to the machine learning programs, an early assumption in the XAI program was that requiring researchers to also incorporate a student model was a bridge too far. But XAI research conducted to date has been a reminder that explanations need to be tailored—somehow—to the knowledge and goals of the user. It is certainly unacceptable to assume that that the user’s understanding of the task is the same as that of the researchers [19]. Furthermore, only XAI research that utilized post-experimental cognitive interviews shows the kind of awareness this research requires.

On the other hand, the neural network learning method addresses perceptual cognition, which symbolic AI, the representational foundation of ITS research, finessed. When images are involved, they are usually presented to the student as text, in terms of already abstracted categories (e.g., the morphology of cultured organism is a “rod”). When the focus is image interpretation itself (e.g., x-ray interpretation, [18]), manually annotated images are presented to the student (e.g., [14]). MR Tutor [22] is a relevant exception in the domain of Magnetic Resonance Imaging (MRI). Experts used a predefined ontology of features to label images, and neural network learning was used to relate patient cases. The resulting “typicality” model enabled the student to view the distribution of disease features across cases, and for the tutoring program to select appropriate problems and examples from the library ([22], pp. 5–8). However, MR Tutor explanations are limited to relating cases, rather than explicating the underlying causal processes that give rise to the observed morphologies—a capability required by specialists for recognizing and discriminating atypical manifestations of a disease.

Studies of radiological expertise revealed an ability to rapidly and automatically recognize “varied normal anatomy” coupled with an ability to describe “abnormal appearance” ([22], p. 3-4). By extension, use of neural network tools for practical applications, a primary objective of XAI research, may require users to have similar capabilities to distinguish discrepant features of interest from normal variation in appearance. Training—and by implication explanation—could be oriented accordingly by presenting normal and abnormal examples and ordering them within a cognitively justified instructional strategy, that is, a pedagogical method.

Pedagogy

Some developers of XAI systems have recognized the need for XAI systems to have a pedagogical foundation (e.g., Raytheon/BBN and Rutgers projects). However, most XAI programs don’t base explanations on an explicit model of the instructional process involving structured methods of interaction, which in turn is based on a theory of learning. By analogy to ITS, an XAI program should incorporate a model for evaluating and instructing proper use of the associated AI program.

Early ITS systems that incorporated pedagogical models include Guidon [6, 9] and Meno-Tutor [15, 26]. The RadTutor [2] for diagnostic interpretation of mammogram images is based on instructional principles (multiplicity, activeness, accommodation and adaptation, and authenticity) and methods (including modelling, coaching, fading of assistance, structured problem solving, and situated learning).

The designers of MR Tutor formulated the following requirements for a computer system to train people in image processing (quoted from [22], p. 4):

- Base the training on a large library of cases representative of [image processing] practice;
- Provide a means of making rapid comparisons between cases by similarity of diagnostically relevant features [the role of the machine learning program];
- Expose the trainee to cases in an order that promotes understanding and retention;
- Help the trainee to make rapid, accurate initial judgements;
- Help the trainee to integrate fragmentary knowledge into more general structural schemata;
- Help the trainee to reflect on experience gained and to integrate general and situated knowledge;
- Be implemented on a personal computer, for use as part of self-study at home or work.

Also applicable to image categorization in general is the idea of a domain-specific description language. In MR Tutor, the Image Description Language (IDL) included *functional descriptors* (e.g., lesion homogeneity, lesion grouping, interior patterning), and *image features* (e.g., visibility, location, shape, size, intensity).

We hypothesize that analogous feature categories and feature descriptions are used by people for interpreting images in general, either formally as standards within a community of practice, or informally by individuals developing their own conscious method for interpreting and classifying images. The use of such feature languages in a variety of domains suggests that comprehending and trusting AI program interpretations, a primary objective of XAI systems, requires an image description language that conforms to the natural

language used in the domain.

Furthermore, instructional research based in cognitive studies suggests that the chain model:

[XAI generates explanations
User comprehends the explanations
User performance improves]

is far too simple—it ignores the active aspect of learning, especially self-explanation. Self-explanation improves learning whether it is prompted or self-motivated [4, 5, 20]. In general, XAI programs do not facilitate self-explanation. Initial instructions given to participants provides explanatory material and may support the self-explanation process; but not all XAI projects provide such instructions. Although some of the projects present examples and tasks that permit displaying boundary conditions (e.g., what the AI gets wrong, false positives), placing the user in a self-explanation mode, XAI methods have not generally exploited the user’s *active efforts* to construct an explanation of the AI system.

CONCLUSION

Some scientific contributions are common to XAI and ITS research. Both seek to promote people’s learning through automated interaction and explanation. Both represent processes as formal models and algorithms in a computer program, in application domains relevant to DoD concerns. Both have found that explanations are more productive when people can respond to them interactively (e.g., by asking follow-up questions), involving theories about when and what kind of explanations facilitate understanding. Researchers in both areas also recognize the need for pilot studies to evaluate the instructional methods and procedures for assessing user understanding.

There have also been contributions of XAI that were not incorporated in the ITS work. Through the use of a symbolic problem-solving model (the embedded expert system), many ITS programs can solve new cases, but for pedagogical effectiveness, most use a curriculum of solved problems curated and organized by specialists (i.e., a “case library”), based on an ontology that has been established within the technical domain (e.g., MR Tutor [22]). It would be advantageous to couple the MR Tutor’s ability to relate cases with the ability of neural network systems to add solved cases to the library.

Another advance is the concern in XAI research with the development of appropriate trust and reliance. Research has demonstrated, for instance, that global explanations alone do not promote trust [19]. ITS research usually focused on teaching people to solve problems themselves, rather than teaching them how to use an AI program that assists them in carrying out complicated technical activities.

In conclusion, the objective of the XAI research program—to develop computational aids to promote practical use of an AI tool, including promoting a user’s understanding of the system’s capabilities and vulnerabilities in practical situations—is inseparable from the objectives of ITS research involving domains of professional expertise, such as medicine, electronics troubleshooting, and engineering. We described the principles of ITS design, in which an explicit pedagogical strategy is based on a cognitive theory of learning in the domain of interest, which is expressed in a model of the subject material. That is, in ITS the design of “explanation systems” is guided by *a well-developed scientific framework*, formalized in process models of problem solving, learning, and communication. We conclude that it will be productive for XAI researchers to view “explanation” as an aspect of an instructional process in which the user is a learner and the program is a tutor, with many of the attendant issues of developing a shared language and understanding of problem-solving methods that ITS research has considered over the past 50 years.

Acknowledgement and Disclaimer

This material is approved for public release. Distribution is unlimited. This material is based on research sponsored by the Air Force Research Lab (AFRL) under agreement number FA8650-17-2-7711. The U.S.

Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of AFRL or the U.S. Government.

REFERENCES

- [1] Anderson, J.R., Corbett, A.T., Koedinger, K.R., and Pelletier, R. (1995). Cognitive tutors: Lessons learned. *The Journal of the Learning Sciences* , 4(2), 167-207.
- [2] Azevedo, R., Lajoie, S., Desaulniers, M., Fleischer, D., and Bret, P. (1997). RadTutor: The theoretical and empirical basis for the design of a mammography interpretation tutor. In B. du Boulay and R. Mizoguchi (Eds.), *Artificial Intelligence in Education: Knowledge and Media in Learning Systems* (pp. 386-393). Amsterdam: IOS Press.
- [3] Brown, J.S., Burton, R.B., and deKleer, J. (1982). Pedagogical, natural language and knowledge engineering techniques. In D. Sleeman and J.S. Brown (Eds.), *SOPHIE I, II and III. Intelligent Tutoring Systems*, (pp. 227–282). New York, Academic Press.
- [4] Chi, M.T., and VanLehn, K.A. (1991). The content of physics self-explanations. *The Journal of the Learning Sciences* ,1 (1), 69–105.
- [5] Chi, M.T., Leeuw, N., Chiu, M.-H., and LaVancher, C. (1994). Eliciting self-explanations improves understanding. *Cognitive Science* , 18(3), 439–477.
- [6] Clancey, W.J. (1979). Dialogue management for rule-based tutorials. *Proceedings of Sixth IJCAI* (pp. 155–161). Tokyo.
- [7] Clancey, W.J. (1985). Heuristic classification. *Artificial Intelligence* , 27:289-350.
- [8] Clancey, W.J. (1986). Qualitative student models. *Annual Review of Computer Science*, 1, 381–450.
- [9] Clancey, W.J. (1987). *Knowledge-based tutoring* . Cambridge, MA: MIT Press.
- [10] Clancey, W.J. (1992). Model construction operators. *Artificial Intelligence* , 53, 1-124.
- [11] Clancey, W.J. (1997). *Situated Cognition: On Human Knowledge and Computer Representations*. New York: Cambridge University Press.
- [12] Clancey, W.J., Bennett, J.S., and Cohen, P.R. (1982). Applications-oriented AI research: Education. In A. Barr and E. A. Feigenbaum (Eds.), *Handbook of artificial intelligence* , Volume 2 (pp. 223–294). Los Altos, CA: William Kaufmann.
- [13] Clancey, W.J. and Soloway, E. (Eds.) (1990). Artificial intelligence and learning environments: Preface. In Special Issue on AI and Learning Environments, *Artificial Intelligence*, Vol. 42.
- [14] Crowley, R.S., Legowski, E., Medvedeva, O., Tseytlin, E., Roh, E., and Jukic, D. (2007). Evaluation of an intelligent tutoring system in pathology: effects of external representation on performance gains, metacognition, and acceptance. *JAMIA* , 14 (2), 182–190.
- [15] Eliot, C. and Woolf, B.P. (1995). An Adaptive Student Centered Curriculum for an Intelligent Training System. *User modeling and user-adapted interaction* 5 , 67–86.
- [16] Hollan, J.D., Hutchins, E.L. and Weitzman, L. (1984). STEAMER: An interactive inspectable simulation-based training system, *AI Magazine* , 5 (2), 15–27.
- [17] Langley, P., Ohlsson, S. and Sage, S. (1984). A machine learning approach to student modeling. Tech. Rept. CMU-RI-TR-84-7, Robotics Institute, Carnegie-Mellon University, Pittsburgh, PA.

- [18] Lesgold, A., Rubinson, H., Feltovich, P., Glaser, R., Klopfer, D., and Wang, Y. (1988). Expertise in a complex skill: Diagnosing x-ray pictures. In M. T. H. Chi, R. Glaser, and M. J. Farr (Eds.), *The nature of expertise* (p. 311–342). Lawrence Erlbaum Associates, Inc.
- [19] Mueller, S. T., Veinott, E. S., Hoffman, R. R., Klein, G., Alam, L., Mamun, T., and Clancey, W. J. (2021). Principles of Explanation in Human-AI Systems, AAAI Virtual Conference, February.
- [20] O’Reilly, T., Symons, S., and MacLachy-Gaudet, H. (1998). A comparison of self-explanation and elaborative interrogation. *Contemporary Educational Psychology* , 23(4), 434–445.
- [21] Schön, D. A. (1987). *Educating the reflective practitioner* . San Francisco: Jossey- Bass.
- [22] Sharples, M., Jeffery, N.P., du Boulay, B., Teather, B.A., Teather, D., and du Boulay, G.H. (2000). Structured computer-based training in the interpretation of neuroradiological images. *International Journal of Medical Informatics* , 60 (3), 263–280.
- [23] Sleeman, D. and Brown, J.S. (Eds.). (1982). *Intelligent tutoring systems* . New York: Academic Press.
- [24] Sottolare, R.A., Perez, R.S., and Skinner, A. (2018). *Assessment of intelligent tutoring systems technologies and opportunities* . NATO Science and Technology Organization, Human Factors and Medicine Task Group 237, Technical Report STO-TR-HFM-237.
- [25] VanLehn, K. (2011). The relative effectiveness of human tutoring, intelligent tutoring systems, and other tutoring systems. *Educational Psychologist*, 46 , 197-221.
- [26] Woolf, B. and McDonald, D. (1984). Context-dependent transitions in tutoring discourse. *Proceedings AAAI-84* (pp. 355–361). Austin, TX.
- [27] Woolf, B.P. (2009). *Building interactive tutors: Student-centered strategies for revolutionizing e-learning* . Burlington, MA: Morgan-Kaufmann.