

# Visual Analytics Perspectives on Interactive and Explainable Machine Learning

Mennatallah El-Assady

University Konstanz



@manunna\_91

[el-assady.com](http://el-assady.com)

# The State of the Art in Integrating Machine Learning into Visual Analytics

A. Endert<sup>1</sup>, W. Ribarsky<sup>2</sup>, C. Turkay<sup>3</sup>, W. Wong<sup>4</sup>, I. Nabney<sup>5</sup>, I. Diaz Blanco<sup>6</sup>, F. Rossi<sup>7</sup>

<sup>1</sup>Georgia Tech, USA

<sup>2</sup>University of North Carolina, Charlotte, USA

<sup>3</sup>City University of London, UK

<sup>4</sup>Middlesex University, UK

<sup>5</sup>Aston University, UK

<sup>6</sup>University of Oviedo, Spain

<sup>7</sup>Paris 1 Panthéon Sorbonne University, Paris

## Abstract

*Visual analytics systems combine machine learning or other analytic techniques with interactive data visualization to promote sensemaking and analytical reasoning. It is through such techniques that people can make sense of large, complex data. While progress has been made, the tactful combination of machine learning and data visualization is still under-explored. This state-of-the-art report presents a summary of the progress that has been made by highlighting and synthesizing select research advances. Further, it presents opportunities and challenges to enhance the synergy between machine learning and visual analytics for impactful future research directions.*

Categories and Subject Descriptors (according to ACM CCS): Human-centered computing - Visualization, Visual analytics

## 1. Introduction

We are in a data-driven era. Increasingly more domains generate and consume data. People have the potential to understand phenomena in more depth using new data analysis techniques. Additionally, new phenomena can be uncovered in domains where data is becoming available. Thus, making sense of data is becoming increasingly important, and this is driving the need for systems that enable people to analyze and understand data.

However, this opportunity to discover also presents challenges. Reasoning about data is becoming more complicated and difficult as data scales and complexities increase. People require powerful tools to draw valid conclusions from data, while maintaining trustworthy and interpretable results.

We claim that visual analytics (VA) and machine learning (ML) have complementing strengths and weaknesses to address these challenges. Visual analytics (VA) is a multi-disciplinary domain that combines data visualization with machine learning (ML) and other automated techniques to create systems that help people make sense of data [TC05, KSF\*08, Kei02, KMSZ06]. Over the years, much work has been done to establish the foundations of this area,

create research advances in select topics, and form a community of researchers to continue to evolve the state of the art.

Currently, VA techniques exist that make use of select ML models or algorithms. However, there are additional techniques that can apply to the broader visual data analysis process. Doing so reveals opportunities for how to couple user tasks and activities with such models. Similarly, opportunities exist to advance ML models based on the cognitive tasks invoked by interactive VA techniques.

This state-of-the-art report briefly summarizes the advances made at the intersection of ML and VA. It describes the extent to which machine learning methods are utilized in visual analytics to date. Further, it illuminates the opportunities within both disciplines that can drive important research directions in the future. Much of the content and inspiration for this paper originated during a Dagstuhl Seminar titled, “Bridging Machine Learning with Information Visualization (15101)” [KMRV15].

arXiv:1702.01226v1 [cs.LG] 4 Feb 2017

# Towards Better Analysis of Machine Learning Models: A Visual Analytics Perspective

Shixia Liu, Xiting Wang, Mengchen Liu, Jun Zhu

Tsinghua University, Beijing, China

## Abstract

Interactive model analysis, the process of understanding, diagnosing, and refining a machine learning model with the help of interactive visualization, is very important for users to efficiently solve real-world artificial intelligence and data mining problems. Dramatic advances in big data analytics has led to a wide variety of interactive model analysis tasks. In this paper, we present a comprehensive analysis and interpretation of this rapidly developing area. Specifically, we classify the relevant work into three categories: understanding, diagnosis, and refinement. Each category is exemplified by recent influential work. Possible future research opportunities are also explored and discussed.

**Keywords:** interactive model analysis, interactive visualization, machine learning, understanding, diagnosis, refinement

## 1. Introduction

Machine learning has been successfully applied to a wide variety of fields ranging from information retrieval, data mining, and speech recognition, to computer graphics, visualization, and human-computer interaction. However, most users often treat a machine learning model as a black box because of its incomprehensible functions and unclear working mechanism [1, 2, 3]. Without a clear understanding of how and why a model works, the development of high-performance models typically relies on a time-consuming trial-and-error pro-

<sup>1</sup>Fully documented templates are available in the elisarticle package on CTAN.

# Visual Analytics in Deep Learning: An Interrogative Survey for the Next Frontiers

Fred Hohman, Member, IEEE, Minsuk Kahng, Member, IEEE, Robert Pienta, Member, IEEE, and Duen Horng Chau, Member, IEEE

**Abstract**—Deep learning has recently seen rapid development and received significant attention due to its state-of-the-art performance on previously-thought hard problems. However, because of the internal complexity and nonlinear structure of deep neural networks, the underlying decision making processes for why these models are achieving such performance are challenging and sometimes mystifying to interpret. As deep learning spreads across domains, it is of paramount importance that we equip users of deep learning with tools for understanding when a model works correctly, when it fails, and ultimately how to improve its performance. Standardized toolkits for building neural networks have helped debug deep learning; visual analytics systems have now been developed to support model explanation, interpretation, debugging, and improvement. We present a survey of the role of visual analytics in deep learning research, which highlights its short yet impactful history and thoroughly summarizes the state-of-the-art using a human-centered interrogative framework, focusing on the *Five W's and How* (Why, Who, What, How, When, and Where). We conclude by highlighting research directions and open research problems. This survey helps researchers and practitioners in both visual analytics and deep learning to quickly learn key aspects of this young and rapidly growing body of research, whose impact spans a diverse range of domains.

**Index Terms**—Deep learning, visual analytics, information visualization, neural networks

## 1 INTRODUCTION

DEEP learning is a specific set of techniques from the broader field of machine learning (ML) that focus on the study and usage of *deep* artificial neural networks to learn structured representations of data. First mentioned as early as the 1940s [1], artificial neural networks have a rich history [2], and have recently seen a dominate and pervasive resurgence [3], [4], [5] in many research domains by producing state-of-the-art results [6], [7] on a number of diverse big data tasks [8], [9]. For example, the premiere machine learning, deep learning, and artificial intelligence (AI) conferences have seen enormous growth in attendance and paper submissions since early 2010s. Furthermore, open-source toolkits and programming libraries for building, training, and evaluating deep neural networks have become more robust and easy to use, democratizing deep learning. As a result, the barrier to developing deep learning models is lower than ever before and deep learning applications are becoming pervasive.

While this technological progress is impressive, it comes with unique and novel challenges. For example, the lack of interpretability and transparency of neural networks, from the learned representations to the underlying decision process, is an important problem to address. Making sense of why a particular model misclassifies test data instances or behaves poorly at times is a challenging task for model developers. Similarly, end-users interacting with an application that relies on deep learning to make critical decisions may question its reliability if no explanation is given by the model, or become baffled if the explanation is convoluted.

• F. Hohman, M. Kahng, R. Pienta, and D. H. Chau are with the College of Computing, Georgia Tech, Atlanta, Georgia 30332, U.S.A.  
E-mail: {fredhohman, kahng, pienta, polo}@gatech.edu

While explaining neural network decisions is important, there are numerous other problems that arise from deep learning, such as AI safety and security (e.g., when using models in applications such as self-driving vehicles), and compromised trust due to bias in models and datasets, just to name a few. These challenges are often compounded, due to the large datasets required to train most deep learning models. As worrisome as these problems are, they will likely become even more widespread as more AI-powered systems are deployed in the world. Therefore, a general sense of model understanding is not only beneficial, but often required to address the aforementioned issues.

Data visualization and visual analytics excel at knowledge communication and insight discovery by using encodings to transform abstract data into meaningful representations. In the seminal work by Zeiler and Fergus [10], a technique called *deconvolutional networks* enabled projection from a model’s learned feature space back to the pixel space. Their technique and results give insight into what types of features deep neural networks are learning at specific layers, and also serve as a debugging tool for improving a model. This work is often credited for popularizing visualization in the machine learning and computer vision communities in recent years, putting a spotlight on it as a powerful tool that helps people understand and improve deep learning models. However, visualization research for neural networks started well before [11], [12], [13]. Over just a handful of years, many different techniques have been introduced to help interpret what neural networks are learning. Many such techniques generate static images, such as attention maps and heatmaps for image classification, indicating which parts of an image are most important to the classification. However, interaction has also been

© 2019 IEEE. This is the author’s version of the article that has been published in IEEE Transactions on Visualization and Computer Graphics. The final version of this record is available at: [10.1109/TVCG.2019.2934629](https://doi.org/10.1109/TVCG.2019.2934629)

# explAIner: A Visual Analytics Framework for Interactive and Explainable Machine Learning

Thilo Spinner, Udo Schlegel, Hanna Schäfer, and Mennatallah El-Assady

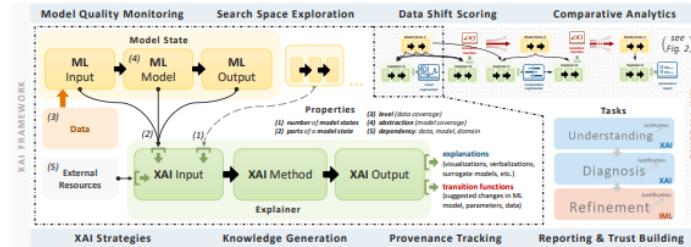


Fig. 1: Close-up view of an *explainer*, the main building-block used to construct an iterative XAI pipeline for the *understanding*, *diagnosis*, and *refinement* of ML models. Explainers have five properties; they take one or more model states as input, applying an XAI method, to output an explanation or a transition function. Global monitoring and steering mechanisms expand the pipeline to the full XAI framework, supporting the overall workflow by guiding, steering, or tracking the explainers during all steps.

**Abstract**—We propose a framework for interactive and explainable machine learning that enables users to (1) understand machine learning models; (2) diagnose model limitations using different explainable AI methods; as well as (3) refine and optimize the models. Our framework combines an iterative XAI pipeline with eight global monitoring and steering mechanisms, including quality monitoring, provenance tracking, model comparison, and trust building. To operationalize the framework, we present explAIner, a visual analytics system for interactive and explainable machine learning that instantiates all phases of the suggested pipeline within the commonly used TensorFlow environment. We performed a user-study with nine participants across different expertise levels to examine their perception of our workflow and to collect suggestions to fill the gap between our system and framework. The evaluation confirms that our tightly integrated system leads to an informed machine learning process while disclosing opportunities for further extensions.

**Index Terms**—Explainable AI, Interactive Machine Learning, Deep Learning, Visual Analytics, Interpretability, Explainability

## 1 INTRODUCTION

Since the first presentation of neural networks in the 1940s [47], we have seen a great increase in works on Artificial Intelligence (AI) and Machine Learning (ML). Especially within the last decade, computational resources have become cheaper and more accessible. This development has led to new state-of-the-art solutions, e.g., Deep Learning (DL), while the increasing availability of tools and libraries has led to a democratization of ML methods in a variety of domains [30]. For example, DL methods outperform traditional algorithms for image processing [56] or natural language processing [82] and can often be applied by domain experts without prior ML expertise [12].

Despite the significant improvement in performance, DL models create novel challenges, due to their nature of being *black-boxes* [78]. For model developers, missing transparency in the decision-making of DL models often leads to a time-consuming trial and error process [81]. Additionally, whenever such decisions concern end-user ap-

plications, e.g., self-driving cars, trust is essential. In critical domains, this trust has to be substantiated either by reliable and unbiased decision outcomes, or convincing rationalization and justifications [66]. The growing prevalence of AI in security-critical domains leads to an ever-increasing demand for explainable and reproducible results.

Several solutions address the problem of missing transparency in black-box models, often referred to as *Explainable Artificial Intelligence* (XAI) [26]. Even though AI algorithms often cannot be directly explained [2], XAI methods aim to provide human-readable, as well as interpretable explanations of the decisions taken by such algorithms. XAI is further driven by newly introduced regulations, such as the *European General Data Protection Regulation* [77], demanding accessible justifications for automated, consumer-facing decisions, prompting businesses to seek reliable XAI solutions. A natural way to obtain human interpretable explanations is through visualizations.

More recent work focuses not only on visual design but also on interactive, mixed-initiative workflows, as provided by Visual Analytics (VA) systems [21]. Also, an exploratory workflow [60] enables a more targeted analysis and design of ML models. Visual analytics further helps in bridging the gap between user knowledge and the insights XAI methods can provide. As AI is affecting a broader range of user groups, ranging from everyday users to model developers, the differing levels of background knowledge in these user groups bring along varying requirements for the explainability.

T. Spinner, U. Schlegel, H. Schäfer, and M. El-Assady are with the University of Konstanz. E-mail: [firstname.lastname@uni-konstanz.de](mailto:firstname.lastname@uni-konstanz.de)  
Manuscript received xx xx, 201x; accepted xx xx, 201x. Date of publication xx xx, 201x; date of current version xx xx, 201x.  
For information on obtaining reprints of this article, please send e-mail to: [reprints@ieee.org](mailto:reprints@ieee.org).  
Digital Object Identifier: [xx.xxxx/TVCG.2019.xxxxxxx/](https://doi.org/10.1109/TVCG.2019.2934629)

## Visual Analytics in Deep Learning

## Interrogative Survey Overview

## §4 WHY

Why would one want to use visualization in deep learning?

Interpretability & Explainability  
Debugging & Improving Models  
Comparing & Selecting Models  
Teaching Deep Learning Concepts



55 WHO

### *Who would use and benefit from visualizing deep learning?*

Model Developers & Builders  
Model Users  
Non-experts

57 HOW

How can we visualize deep learning data, features, and relationships?

- Node-link Diagrams for Network Architecture
- Dimensionality Reduction & Scatter Plots
- Line Charts for Temporal Metrics
- Instance-based Analysis & Exploration
- Interactive Experimentation
- Algorithms for Attribution & Feature Visualization

59 WHERE

### Where has deep learning visualization been used?

Application Domains & Models  
A Vibrant Research Community

	WHY	WHO	WHAT	HOW	WHEN	WHERE
Work	4.1	4.2	4.3	4.4	5.1	5.2
Abadi, et al., 2016 [27]						arXiv
Bau, et al., 2017 [28]						CVPR
Bilal, et al., 2017 [29]						TVCG
Bojarski, et al., 2016 [30]						arXiv
Bruckner, 2014 [31]						MS Thesis
Carter, et al., 2016 [32]						Distill
Cashman, et al., 2017 [33]						VADL
Chae, et al., 2017 [34]						VADL
Chung, et al., 2016 [35]						FILM
Goyal, et al., 2016 [36]						arXiv
Harley, 2015 [37]						ISVC
Hohman, et al., 2017 [38]						CHI
Kahng, et al., 2018 [39]						TVCG
Karpathy, et al., 2015 [40]						arXiv
Li, et al., 2015 [41]						arXiv
Liu, et al., 2017 [14]						TVCG
Liu, et al., 2018 [42]						TVCG
Ming, et al., 2017 [43]						VAST
Norton & Qi, 2017 [44]						VizSec
Olah, 2014 [45]						Web
Olah, et al., 2018 [46]						Distill
Pezzotti, et al., 2017 [47]						TVCG
Rauber, et al., 2017 [48]						TVCG
Robinson, et al., 2017 [49]						GeoHum.
Rong & Adar, 2016 [50]						NIPS WS.
Smilkov, et al., 2016 [51]						ICML VIS
Smilkov, et al., 2017 [16]						ICML VIS
Strobelt, et al., 2018 [52]						TVCG
Tzeng & Ma, 2005 [13]						VIS
Wang, et al., 2018 [53]						TVCG
Webster, et al., 2017 [54]						Web
Wongsuphasawat, et al., 2018 [15]						TVCG
Yosinski, et al., 2015 [55]						ICML DL
Zahavy, et al., 2016 [56]						ICML
Zeiler, et al., 2014 [10]						ECCV
Zeng, et al., 2017 [57]						VADL
Zhong, et al., 2017 [58]						ICML VIS
Zhu, et al., 2016 [59]						ECCV

## The Newcomer.



**Model Novice**

The model novice is the 'new one' in the machine-learning class. His goal is to learn about concepts of machine learning models; he wants to understand the building blocks of the model as well as its general working. Learning resources are essential to him, either by example or by textual, visual, or external resources.

## The Operator.

He is the 'user' among the users, i.e., he uses existing machine learning models to solve specific tasks. For example, this could be a domain expert - let's say a biologist - who needs to classify protein structures. To decide on a model, he wants to compare architectures, understand their underlying working, and verify his decision by executing XAI methods on some data samples.



**Model User**



**Model Developer**

## The Expert.

The model developer is an expert on machine learning. He develops models from scratch, refines existing models, and optimizes parameters to improve the model's performance. He is interested in the architecture of the model, including in-depth information, such as layer-sizes, initializers, and activation functions. To debug the model, explanations on all abstraction levels are relevant. His insights might lead to a model update, covering the full development and refinement process.

# The Process of Interactive and Explainable Machine Learning

## Towards Better Analysis of Machine Learning Models: A Visual Analytics Perspective

Shixia Liu, Xiting Wang, Mengchen Liu, Jun Zhu  
Tsinghua University, Beijing, China

### Abstract

Interactive model analysis, the process of understanding, diagnosing, and refining a machine learning model with the help of interactive visualization, is very important for users to efficiently solve real-world artificial intelligence and data mining problems. Dramatic advances in big data analytics has led to a wide variety of interactive model analysis tasks. In this paper, we present a comprehensive analysis and interpretation of this rapidly developing area. Specifically, we classify the relevant work into three categories: understanding, diagnosis, and refinement. Each category is exemplified by recent influential work. Possible future research opportunities are also explored and discussed.

**Keywords:** interactive model analysis, interactive visualization, machine learning, understanding, diagnosis, refinement

## explAIner: A Visual Analytics Framework for Interactive and Explainable Machine Learning

Thilo Spinner, Udo Schlegel, Hanna Schäfer, and Mennatallah El-Assady

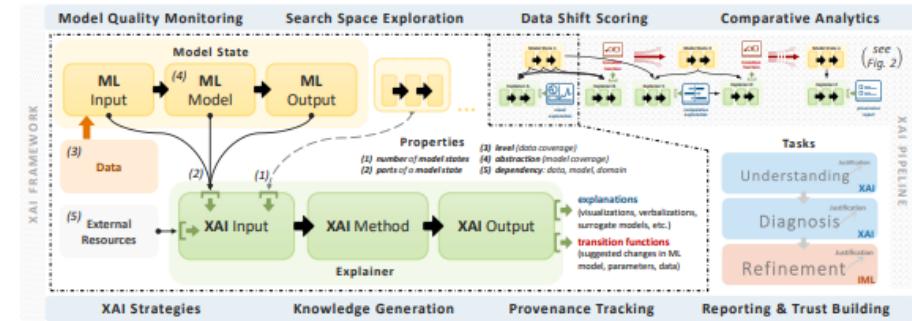


Fig. 1: Close-up view of an *explainer*, the main building-block used to construct an iterative XAI pipeline for the *understanding*, *diagnosis*, and *refinement* of ML models. Explainers have five properties; they take one or more model states as input, applying an XAI method, to output an explanation or a transition function. Global monitoring and steering mechanisms expand the pipeline to the full XAI framework, supporting the overall workflow by guiding, steering, or tracking the explainers during all steps.

**Abstract**— We propose a framework for interactive and explainable machine learning that enables users to (1) understand machine learning models; (2) diagnose model limitations using different explainable AI methods; as well as (3) refine and optimize the models. Our framework combines an iterative XAI pipeline with eight global monitoring and steering mechanisms, including quality monitoring, provenance tracking, model comparison, and trust building. To operationalize the framework, we present explAIner, a visual analytics system for interactive and explainable machine learning that instantiates all phases of the suggested pipeline within the commonly used TensorBoard environment. We performed a user-study with nine participants across different expertise levels to examine their perception of our workflow and to collect suggestions to fill the gap between our system and framework. The evaluation confirms that our tightly integrated system leads to an informed machine learning process while disclosing opportunities for further extensions.

**Index Terms**— Explainable AI, Interactive Machine Learning, Deep Learning, Visual Analytics, Interpretability, Explainability

Understanding  
Justification

Diagnosis  
Justification

Refinement  
Justification

# Understanding

# Do Convolutional Neural Networks Learn Class Hierarchy?

Bilal Alsallakh, Amin Jourabloo, Mao Ye, Xiaoming Liu, Liu Ren

# Understanding Model Decisions

## ThreadReconstructor: Modeling Reply-Chains to Untangle Conversational Text through Visual Analytics

Mennatallah El-Assady<sup>1,2</sup>, Rita Sevastjanova<sup>1</sup>

<sup>1</sup>University of Konstanz  
<sup>2</sup>University of Ontario Institute of Technology

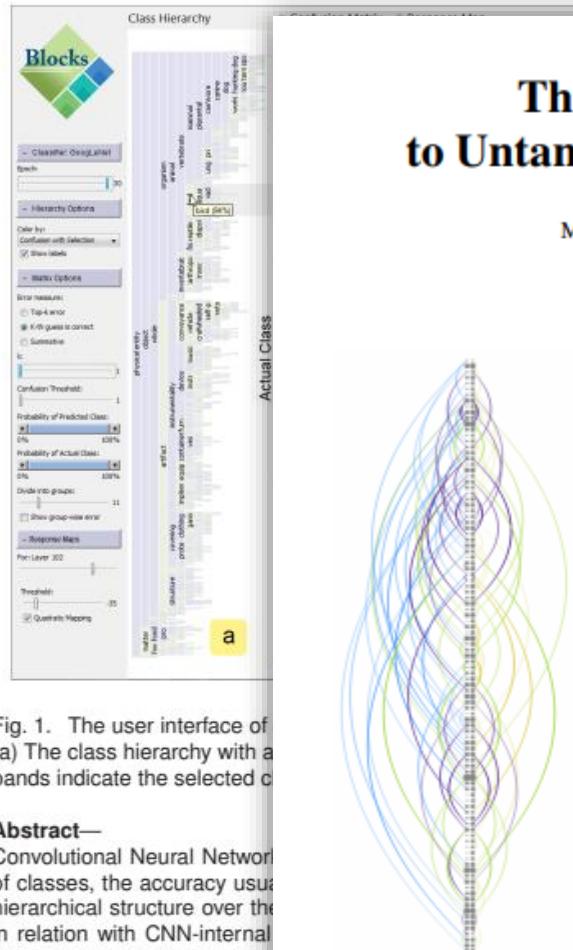


Fig. 1. The user interface of ThreadReconstructor. (a) The class hierarchy with a legend. (b) Reply-Relation View. (c) Thematic-Forest View showing each message as a node and reply-chains as arcs. (d) Control Panel.

### Abstract

Convolutional Neural Networks (CNNs) have shown great success in learning class hierarchies of classes, the accuracy usually being higher than that of other models. The learned hierarchical structure over the classes is usually consistent with the internal structure of the CNNs. Furthermore, the learned hierarchical structure over the classes further dictates the learned features. For example, the learned features can separate high-level groups of classes. This is achieved by training the CNN for more epochs to develop specific features that are key to significant improvements in the accuracy of the model. We further demonstrate that the learned features are key to significant improvements in the accuracy of the model, compared to a random-forest model trained on 13 features (right arcs). More details are provided in the following sections.

**Index Terms**—Convolutional

### Abstract

We present *ThreadReconstructor*, a visual analytics approach for detecting reply-chains in conversational text, e.g., in political debates and forums. Our work is motivated by the need for better understanding of the behavior of machine learning models in massive online conversations and verbatim text transcripts. We propose a visual analytics approach that generates a basic structure that is enriched by user-defined rules to untangle reply-chains in conversational text.

## RuleMatrix: Visualizing and Understanding Classifiers with Rules

Yao Ming, Huamin Qu, *Member, IEEE*, and Enrico Bertini, *Member, IEEE*

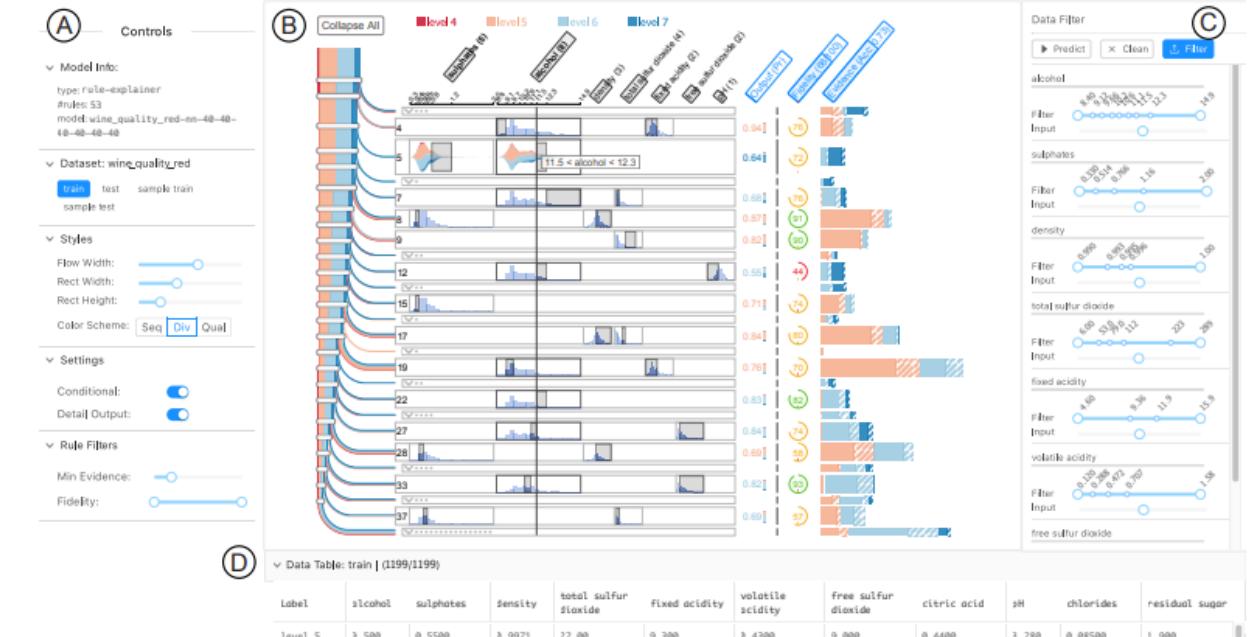


Fig. 1. Understanding the behavior of a trained neural network using the explanatory visual interface of our proposed technique. The user uses the control panel (A) to specify the detail information to visualize (e.g., level of detail, rule filters). The rule-based explanatory representation is visualized as a matrix (B), where each row represents a rule, and each column is a feature used in the rules. The user can also filter the data or use a customized input in the data filter (C) and navigate the filtered dataset in the data table (D).

**Abstract**—With the growing adoption of machine learning techniques, there is a surge of research interest towards making machine learning systems more transparent and interpretable. Various visualizations have been developed to help model developers understand, diagnose, and refine machine learning models. However, a large number of potential but neglected users are the domain experts with little knowledge of machine learning but are expected to work with machine learning systems. In this paper, we present an interactive visualization technique for rule-based classifiers, with the goal of making machine learning models more transparent and interpretable to all the users.

## Visualizing the Hidden Activity of Artificial Neural Networks

Paulo E. Rauber, Samuel G. Fadel, Alexandre X. Falcão, and Alexandru C. Telea

**Abstract**—In machine learning from examples, Artificial neural networks are typically used as black-boxes to represent observations. In this paper, we propose a new way of representing observations using traditional image classification methods. We show that this approach can be used to design better classifiers. For instance, we can use this approach to design better representations, and the paper concludes with some experimental results.

## 1 INTRODUCTION

In machine learning, advances in building and training (deep) arti-  
allowed these models to achieve-  
plications related to pattern recog-  
training ANNs is generally time-  
expertise [5].

In data visualization, dimensionality reduction is often used to compute *projections* of high-dimensional data in lower-dimensional spaces. These projections are used to preserve the data *structure*. This is often done by using a dimensionality reduction technique such as PCA. When depicted by scatterplots, projections can help to reveal the underlying structure of the original data. Compared to other dimensionality reduction techniques, such as t-SNE or MDS, PCA is a linear method that projects data onto a lower-dimensional space. The resulting projections are orthogonal to each other, which makes it easier to interpret the results. For example, if a dataset has three dimensions, PCA will find the first two principal components that explain the most variance in the data. These components are orthogonal to each other, which means that they are independent of each other. This makes it easier to interpret the results of PCA, as it is easier to understand the relationship between the different dimensions of the data. In addition, PCA is a linear method, which means that it is faster to compute than other dimensionality reduction techniques. This makes it a popular choice for large datasets, as it is able to handle millions of observations in a reasonable amount of time. Overall, PCA is a powerful technique for dimensionality reduction that can help to reveal the underlying structure of high-dimensional data.

In this paper, we demonstrate induction techniques to provide insights. Although we focus on multilayer

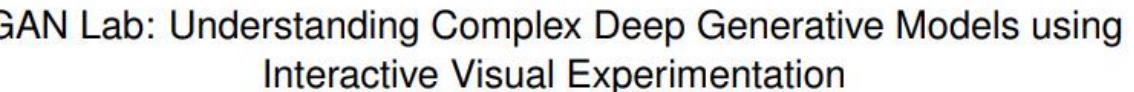


Figure 1: The interface of RNNVis. The control panel (A) shows parameters

## Understanding Hidden Memories of Recurrent Neural Networks

Yao Ming \* Shaozu Cao \* Ruixiang Zhang \* Zhen Li \* Yuanzhe Chen \*  
Yangqiu Song, Member, IEEE \*

Hong Kong University of Scienc



Minsuk Kahng, Nikhil Thorat, Duen Horng (Polo) Chau, Fernanda B. Viégas, and Martin Wattenberg

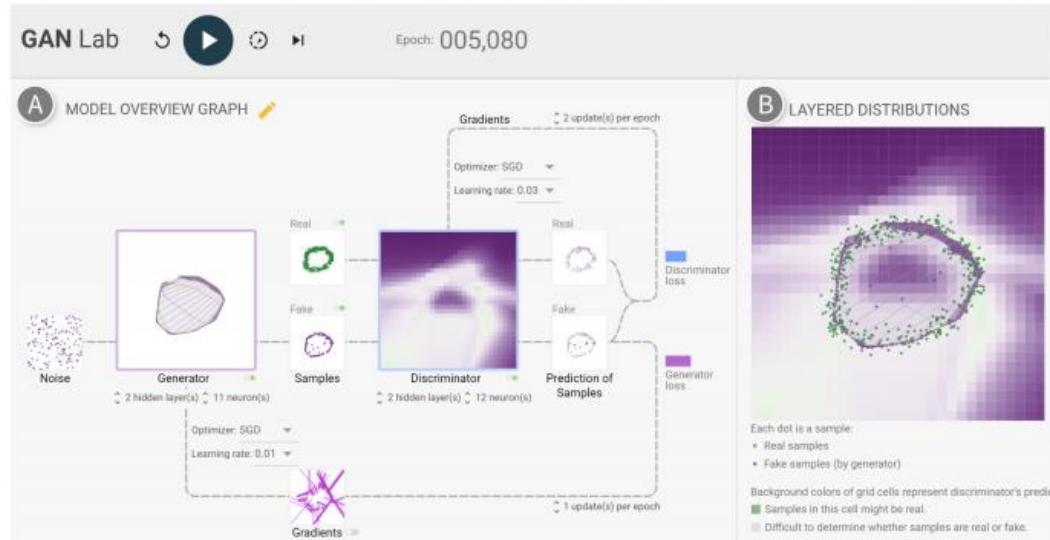


Fig. 1. With GAN Lab, users can interactively train Generative Adversarial Networks (GANs), and visually examine the model training process. In this example, a user has successfully used GAN Lab to train a GAN that generates 2D data points whose challenging distribution resembles a ring. **A.** The *model overview graph* summarizes a GAN model's structure as a graph, with nodes representing the *generator* and *discriminator* submodels, and the data that flow through the graph (e.g., fake samples produced by the generator). **B.** The *layered distributions* view helps users interpret the interplay between submodels through user-selected layers, such as the discriminator's classification heatmap, *real samples*, and *fake samples*, produced by the generator.

**Abstract**—Recent success in deep learning has generated immense interest among practitioners and students, inspiring many to learn about this new technology. While visual and interactive approaches have been successfully developed to help people more easily learn deep learning, most existing tools focus on simpler models. In this work, we present GAN Lab, the first interactive visualization tool designed for non-experts to learn and experiment with *Generative Adversarial Networks (GANs)*, a popular class of complex deep learning models. With GAN Lab, users can interactively train generative models and visualize the dynamic training process's intermediate results. GAN Lab tightly integrates an *model overview graph* that summarizes GAN's structure, and a *layered distributions* view that helps users interpret the interplay between submodels. GAN Lab introduces new interactive experimentation features for learning complex deep learning models, such as *step-by-step* training at multiple levels of abstraction for understanding intricate training dynamics. Implemented using *TensorFlow* is, GAN Lab is accessible to anyone via modern web browsers, without the need for

# Understanding Model Behavior

## LSTMVis: A Tool for Visual Analysis of Hidden State Dynamics in Recurrent Neural Networks

Hendrik Strobelt, Sebastian Gehrmann, Hanspeter Pfister, and Alexander M. Rush  
– Harvard School of Engineering and Applied Sciences –

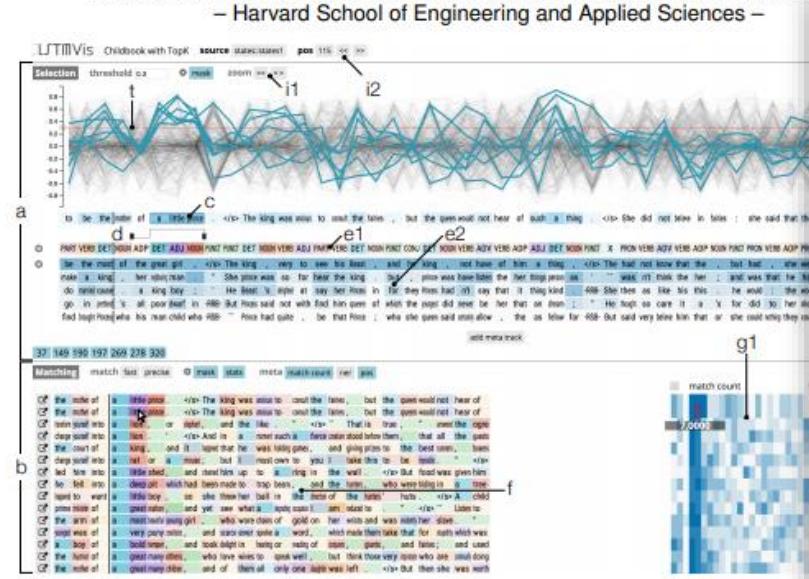


Fig. 1. The LSTMVis user interface. The user interactively *selects* a range of text specifying a hypothesis. Select View (a). This range is then used to *match* similar hidden state patterns displayed in the Match View (b) by specifying a start-stop range in the text (c) and an activation threshold (t) which leads to a selection of *hi*. The start-stop range can be further constrained using the pattern plot (d). The meta-tracks below depict *ext* position like POS (e1) or the top K predictions (e2). The tool can then *match* this selection with similar hidden state set of varying lengths (f), providing insight into the representations learned by the model. The match view displays user-defined meta-data encoded as heatmaps (g1,g2). The color of one heatmap (g2) can be mapped (h) to the color of the other heatmap (g1). This allows the user to see patterns that lead to further refinement of the selection hypothesis. Navigation aids provide

**Abstract**— Recurrent neural networks, and in particular long short-term memory (LSTM) networks, are a remarkable sequence modeling that learn a dense black-box hidden representation of their sequential input. Research on understanding these models have studied the changes in hidden state representations over time and not patterns but also significant noise. In this work, we present LSTMVis, a visual analysis tool for recurrent neural networks that focus on understanding these hidden state dynamics. The tool allows users to select a hypothesis input range, to match these states to similar patterns in a large data set, and to align these results with their domain. We show several use cases of the tool for analyzing specific hidden state properties on data, phrase structure, and chord progressions, and demonstrate how the tool can be used to isolate patterns for further analysis. We characterize the domain, the different stakeholders, and their goals and tasks. Long-term usage data after revealed great interest in the machine learning community.

## SEQ2SEQ-VIS : A Visual Debugging Tool for Sequence-to-Sequence Models

Hendrik Strobelt\*, Sebastian Gehrmann\*, Michael Behrisch, Adam Perer, Hanspeter Pfister, Alexander M. Rush

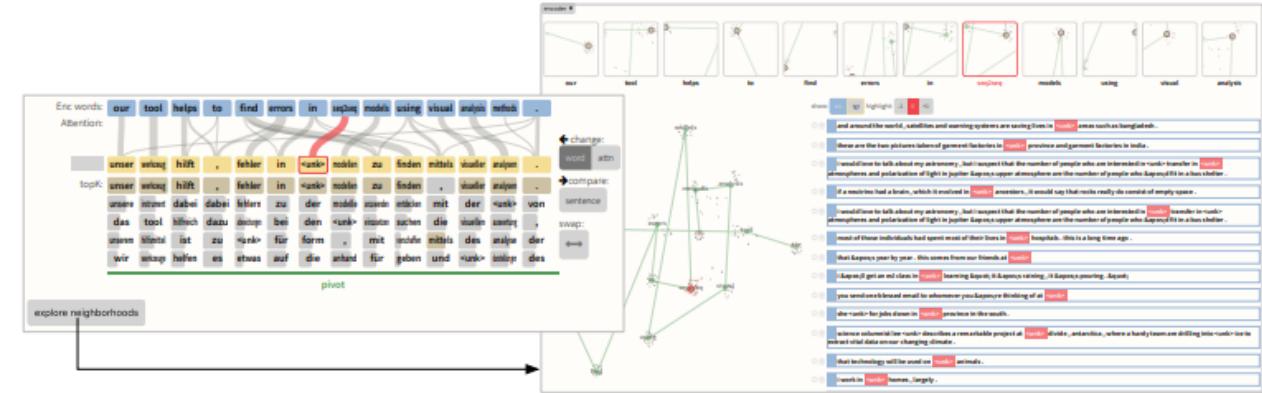


Fig. 1. Example of Seq2Seq-Vis. In the translation view (left), the source sentence “our tool helps to find errors in seq2seq models using visual analysis methods.” is translated into a German sentence. The word “seq2seq” has correct attention between encoder and decoder (red highlight) but is not part of the language dictionary. When investigating the encoder neighborhoods (right), the user sees that “seq2seq” is close to other unknown words (“unk”). The buttons enable user interactions for deeper analysis.

**Abstract**— Neural sequence-to-sequence models have proven to be accurate and robust for many sequence prediction tasks, and have become the standard approach for automatic translation of text. The models work with a five-stage blackbox pipeline that begins with encoding a source sequence to a vector space and then decoding out to a new target sequence. This process is now standard, but like many deep learning methods remains quite difficult to understand or debug. In this work, we present a visual analysis tool that allows interaction and “what if”-style exploration of trained sequence-to-sequence models through each stage of the translation process. The aim is to identify which patterns have been learned, to detect model errors, and to probe the model with counterfactual scenarios. We demonstrate the utility of our tool through several real-world sequence-to-sequence use cases on large-scale models.

# Diagnosis

# A Workflow for Visual Diagnostics of Binary Classifiers using Instance-Level Explanations

Josua Krause\*  
NYU Tandon  
School of Engineering

Aritra Dasgupta†  
Pacific Northwest National Laboratory

Jordan Swartz‡

Yindalon Aphinyanaphongs§

Enrico Bertini¶

## Manifold: A Model-Agnostic Framework for Interpretation and Diagnosis of Machine Learning Models

### ABSTRACT

Human-in-the-loop data analysis applies transparency in machine learning models and trust their decisions. To this end, we propose a workflow to help data scientists analyze, diagnose, and understand the decisions of machine learning models. The approach leverages “instance-level” explanations that explain single instances to build a set of visual representations for model investigation. The workflow is based on two main components and steps: one based on aggregated data distributions across correct / incorrect predictions to understand which features drive decisions; and one based on raw data, to identify potential root causes for the observed patterns. This workflow is the result of a long-term collaboration with a team of data scientists and healthcare professionals who used it to interpret the decisions of machine learning models they developed. This collaboration demonstrates that the workflow is useful for experts to derive useful knowledge about the model they describe, thus experts can generate feedback that the model can be improved.

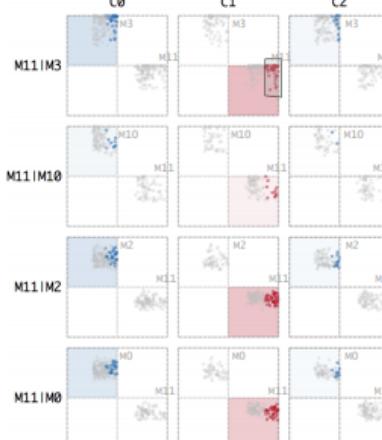
**Keywords:** Machine Learning, Interpretability

### 1 INTRODUCTION

In this paper we propose an interactive interface to help data scientists and domain experts analyze and interpret the decisions of machine learning models.

Jiawei Zhang, Yang Wang, Piero Molino

Between Model Comparison  
An overview of model performance over classes.



\* A larger coordinate indicates higher prediction confidence

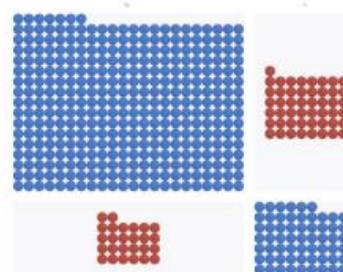
(1)

Fig. 1. Manifold consists of two interactive dialogs: a model comparison dialog (a) that allows users to compare pairs using a small multiple design, and a local feature interpretation dialog (b) that allows users to explore subsets (c) and provides a similarity measure (b) of feature contributions. (d) shows the top 10 most discriminative features among different subsets, i.e., so-called “Manifolds”.

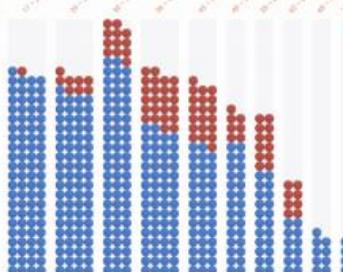
**Abstract**— Interpretation and diagnosis of machine learning models is a challenging task that requires new approaches. We present Manifold, a framework that facilitates the interpretation and comparison of machine learning models in a more transparent way. Manifold focuses on visualizing the internal logic of a specific model type (i.e., regression) in a scenario where different model types are integrated. To this end, Manifold allows users to access the internal logic of the model and solely observes the model’s internal logic and probability distribution. We describe the workflow of Manifold, which is commonly involved in the model development and diagnosis (verification). The visual components supporting these tasks include a local feature interpretation and a customizable tabular view that reveals feature contributions and a classification rule, as well as a small multiple dialog for comparing different models.

## The What-If Tool: Interactive Probing of Machine Learning Models

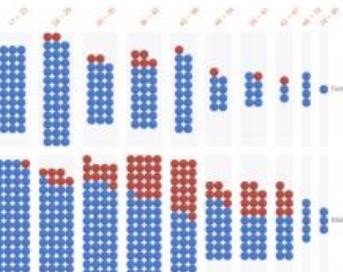
James Wexler, Mahima Pushkarna, Tolga Bolukbasi, Martin Wattenberg, Fernanda Viégas, and Jimbo Wilson



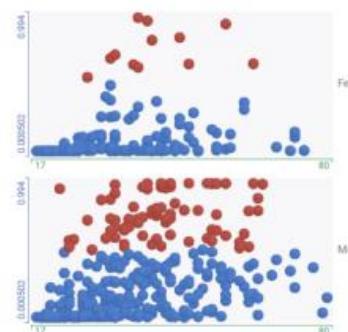
(a) Confusion matrix of a single binary classification model, colored by prediction correctness



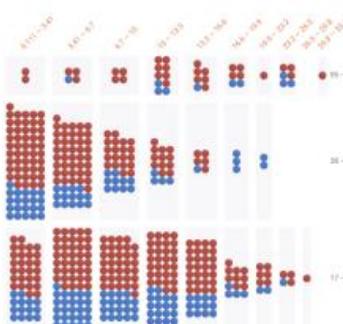
(b) Histogram of age, colored by classification



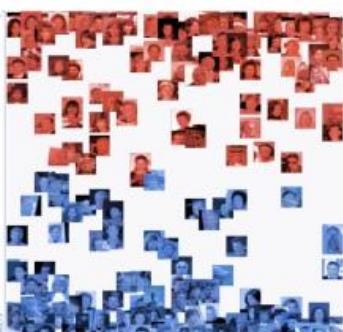
(c) Two-dimensional histogram of age and sex, colored by classification



(d) Small multiples by sex. Each scatterplot shows age vs positive classification score, colored by classification



(e) Histograms of performance in a regression model that predicts age, faceted into 3 age buckets



(f) Using images as thumbnails for image datasets

Fig. 1: Analyzing two different data sets in the What-If Tool: US Census Income dataset from UCI (a) to (e), and CelebA dataset (f).

Diagnosis  
Model Outputs

## Analyzing the Noise Robustness of Deep Neural Networks

Mengchen Liu\*, Shixia Liu\*, Hang Su<sup>†</sup>, Kelei Cao\*, Jun Zhu<sup>†</sup>

<sup>†</sup>School of Software, TNList Lab, State Key Lab for Intell. Tech. Sys., Tsinghua University, Beijing 100084, China  
<sup>†</sup>Dept. of Comp. Sci.Tech., TNList Lab, State Key Lab for Intell. Tech. Sys., CBICR Center, Tsinghua University, Beijing 100084, China

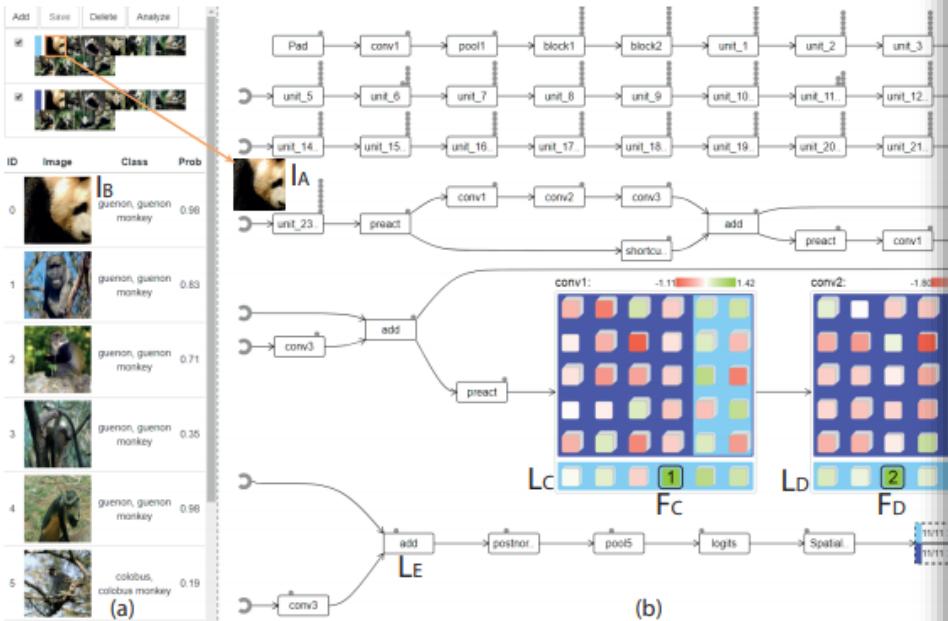
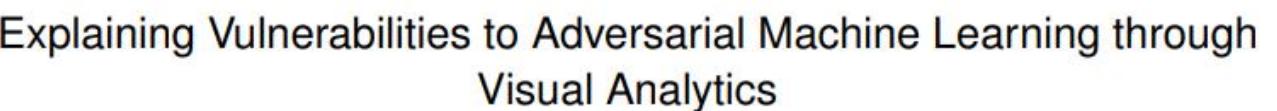


Figure 1: Explaining the misclassification of adversarial panda images. The root cause is that the adversarial examples ( $F_C$ ), which leads to the failure of detecting a panda's face ( $F_D$ ). As a result, the images are misclassified: (a) input images; (b) datapath visualization at the layer and feature map levels.



Yuxin Ma, Tiankai Xie, Jundong Li, Ross Maciejewski, Senior Member, IEEE

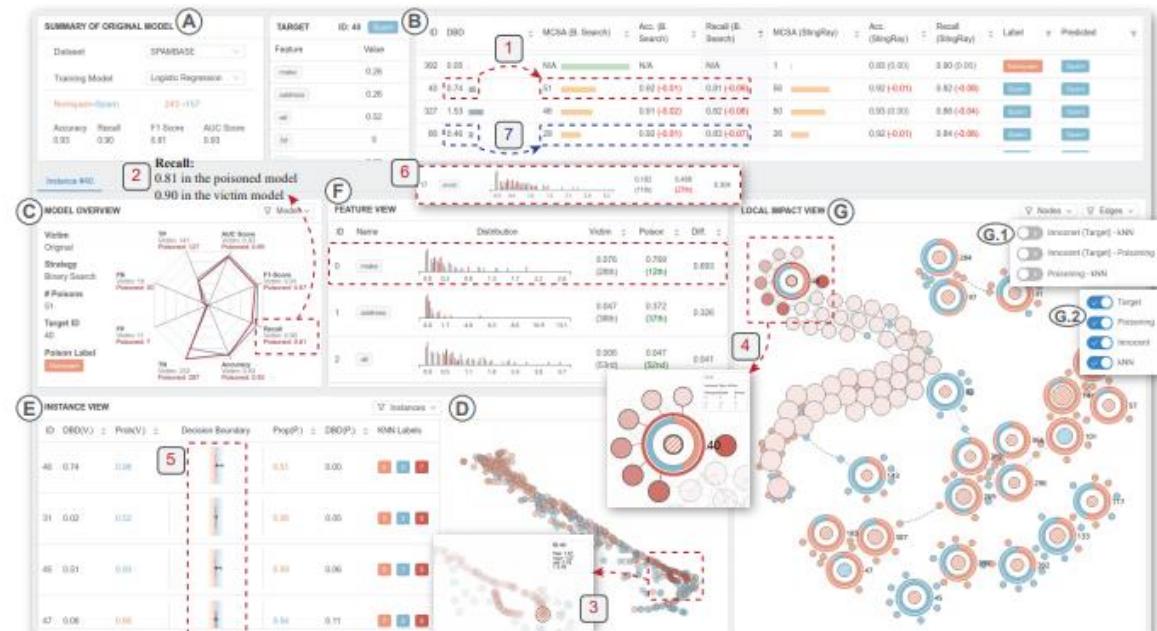


Fig. 1. Reliability attack on spam filters. (1) Poisoning instance #40 has the largest impact on the recall value, which is (2) also depicted in the model overview. (3) There is heavy overlap among instances in the two classes as well the poisoning instances. (4) Instance #40 has been successfully attacked causing a number of innocent instances to have their labels flipped. (5) The flipped instances are very close to the decision boundary. (6) On the feature of words "will" and "email", the variances of poisoning instances are large. (7) A sub-optimal target (instance #80) has less impact on the recall value, but the cost of insertions is 40% lower than that of instance #40.

**Abstract**— Machine learning models are currently being deployed in a variety of real-world applications where model predictions are used to make decisions about healthcare, bank loans, and numerous other critical tasks. As the deployment of artificial intelligence technologies becomes ubiquitous, it is unsurprising that adversaries have begun developing methods to manipulate machine learning models to their advantage. While the visual analytics community has developed methods for opening the black box of machine learning models, little work has focused on helping the user understand their model vulnerabilities in the context of adversarial attacks. In this paper, we present a visual analytics framework for explaining and exploring model vulnerabilities to adversarial attacks. Our framework

# Diagnosis Activation

## SUMMIT: Scaling Deep Learning Interpretability by Visualizing Activation and Attribution Summarizations

Fred Hohman, Haekyu Pa

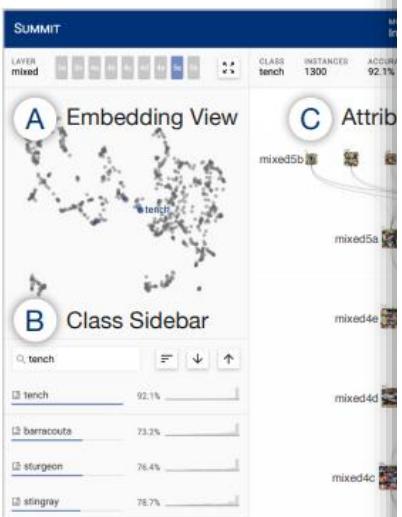


Fig. 1. With SUMMIT, users can scalably summarize network detects and *how* they are related. In this example, the "tench" prediction is dependent on an intermediate "t" like "scales," "person," and "fish". (A) **Embedding View** (B) **Class Sidebar** enables users to search, sort, and filter highly activated neurons as vertices ("scales," "fish")

**Abstract**—Deep learning is increasingly used in deep learning. However, explaining predictions remains a fundamental challenge. Existing approaches for explaining predictions for single images or neurons, such as LIME, are not suitable for explaining predictions for millions of images, such as neural networks. SUMMIT is a system that systematically summarizes and visualizes what features are most important for a model's predictions. SUMMIT introduces two new scalable summarization methods: (1) *neuron-influence aggregation* identifies relationships between neurons and the novel *attribution graph* that reveals and summarizes how neurons interact with each other to produce outcomes. SUMMIT scales to large data, such as neural network models, and provides visualization and dataset examples to help users understand the model's predictions. We present neural network exploration scenarios and compare SUMMIT to other large-scale image classifier's learned representation and visualization systems. SUMMIT runs in modern web browsers and is open-sourced.

**Index Terms**—Deep learning interpretability, visual

## DQNViz: A Visual Analytics Approach to Understand Deep Q-Networks

Junpeng Wang, Liang Gou, Han-Wei S

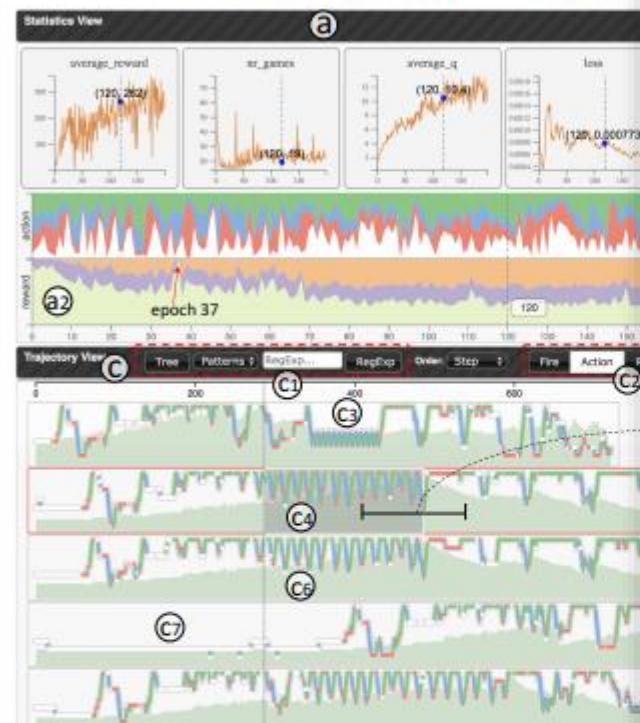


Fig. 1. DQNViz: (a) the **Statistics view** presents the overall training progress and performance of the DQN agent; (b) the **Trajectory view** shows epoch-level statistics with pie charts and stacked bar charts; (c) the **Segment view** shows the activation patterns of the DQN agent in different episodes; (d) the **Segment view** shows the activation patterns of the DQN agent in different episodes.

**Abstract**—Deep Q-Network (DQN), as one type of deep reinforcement learning, has achieved superhuman performance in playing games like Go and chess. Despite the sophisticated behaviors of the DQN agent remain to be challenging, the large number of experiences dynamically generated by the agent.

## ACTIVIS: Visual Exploration of Industry-Scale Deep Neural Network Models

Minsuk Kahng, Pierre Y. Andrews, Aditya Kalro, and Duen Horng (Polo) Chau

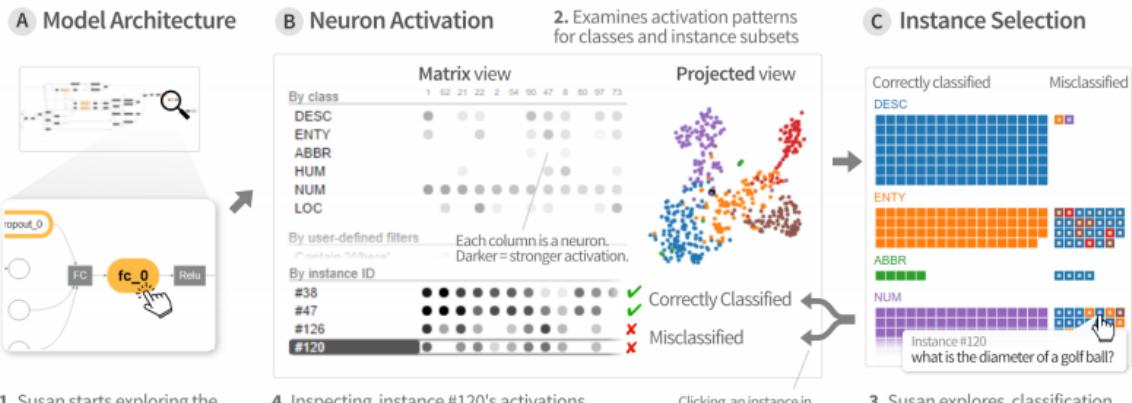


Fig. 1. ACTIVIS integrates several coordinated views to support exploration of complex deep neural network models, at both instance- and subset-level. 1. Our user Susan starts exploring the model overview. She selects a data node (yellow). 2. The **neuron activation matrix view** shows the activations for instances and instance subsets; the **projected view** displays the 2-D projection of instance activations. 3. From the **instance selection panel** (at C), she explores individual instances and their classification results. 4. Adding instances to the matrix view enables comparison of activation patterns across instances, subsets, and classes, revealing causes for misclassification.

**Abstract**—While deep learning models have achieved state-of-the-art accuracies for many prediction tasks, understanding these models remains a challenge. Despite the recent interest in developing visual tools to help users interpret deep learning models, the complexity and wide variety of models deployed in industry, and the large-scale datasets that they used, pose unique design challenges that are inadequately addressed by existing work. Through participatory design sessions with over 15 researchers and engineers at Facebook, we have developed, deployed, and iteratively improved ACTIVIS, an interactive visualization system for interpreting large-scale deep learning models and results. By tightly integrating multiple coordinated views, such as a *computation graph* overview

## Analyzing the Training Processes of Deep Generative Models

Mengchen Liu, Jiaxin Shi, Kelei Cao, Jun Zhu, Shixia Liu

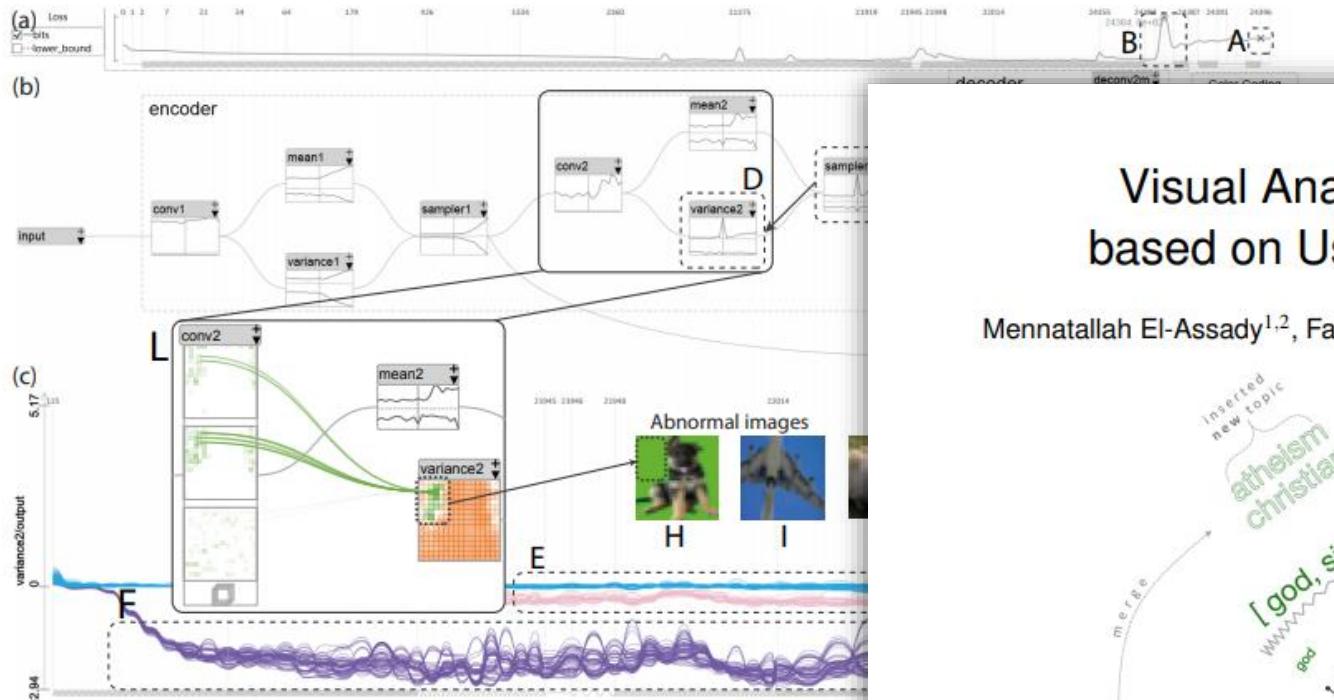


Fig. 1. DGMTracker, a visual analytics tool that helps experts understand and diagnosis models (DGMs): (a) the loss changes; (b) the data flow visualization to illustrate how data neurons influence the output of the neuron of interest; (c) visualization of the training dynamics.

**Abstract**— Among the many types of deep models, deep generative models (DGMs) play a significant role in unsupervised and semi-supervised learning. However, training DGMs requires more skill and time than training other types of deep models such as convolutional neural networks. In this paper, we propose a novel approach for better understanding and diagnosing the training process of a DGM. To achieve this, we first extract a large amount of time series data that represents training data. A blue-noise polyline sampling scheme is then introduced to select time series samples, which helps to reduce visual clutter. To further investigate the root cause of a failed training process, we propose a novel approach that identifies which neurons in the network contribute to the output of the neuron causing the training failure. This approach can be used to understand and diagnose the training process of a DGM. We also demonstrate how our approach can be used to analyze other types of deep models, such as convolutional neural networks.

# Diagnosis Processes

# Visual Analytics for Topic Model Optimization based on User-Steerable Speculative Execution

Mennatallah El-Assady<sup>1,2</sup>, Fabian Sperrle<sup>1</sup>, Oliver Deussen<sup>1</sup>, Daniel Keim<sup>1</sup>, and Christopher Collins<sup>2</sup>

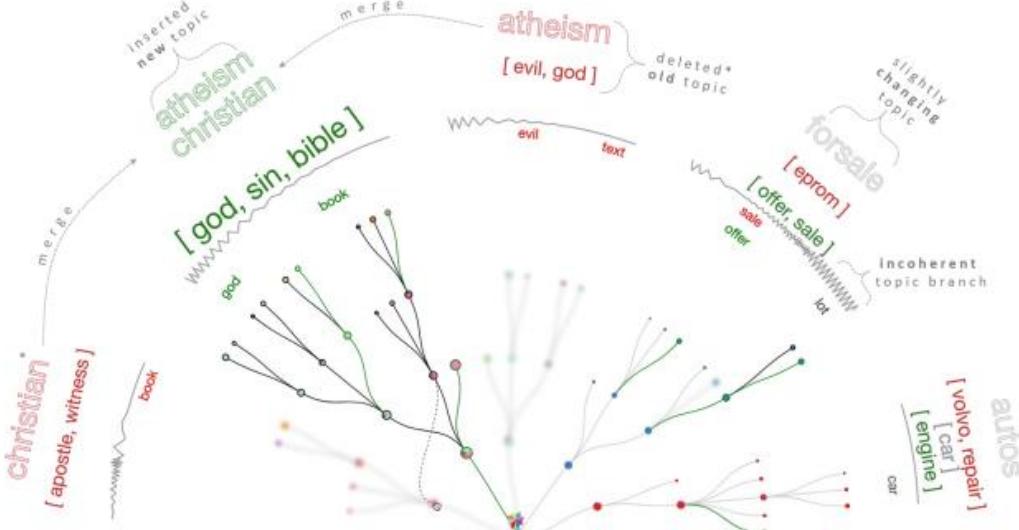


Fig. 1: The *Tree-Speculation View* is used to compare two topic models and shows the differences. **Deleted branches** are blurred, while **moved**, **newly added** and **removed** nodes and keywords are highlighted. To efficiently guide users towards perceivable model quality improvements, our system automatically proposes optimizations like the merge of two topics depicted here. By visualizing model uncertainties and low quality topics, we foster trust in the model and empower users to directly address these shortcomings.

**Abstract**— To effectively assess the potential consequences of human interventions in model-driven analytics systems, we establish the concept of *speculative execution* as a visual analytics paradigm for creating user-steerable preview mechanisms. This paper presents an explainable, mixed-initiative topic modeling framework that integrates speculative execution into the algorithmic decision-making process. Our approach visualizes the model-space of our novel incremental hierarchical topic modeling algorithm, unveiling its inner-workings. We support the active incorporation of the user's domain knowledge in every step through explicit model manipulation interactions. In addition, users can initialize the model with expected topic seeds, the backbone priors. For a more targeted optimization, the modeling process automatically triggers a speculative execution of various optimization strategies, and requests feedback whenever

### FAIRVIS: Visual Analytics for Discovering Intersectional Bias in Machine Learning

Ángel Alexander Cabrera  
Jamie Morgenstern

Will Epperson

Fred Hohman  
Duen Horng (Polo) Chau\*

Minsuk Kahng

Georgia Institute of Technology

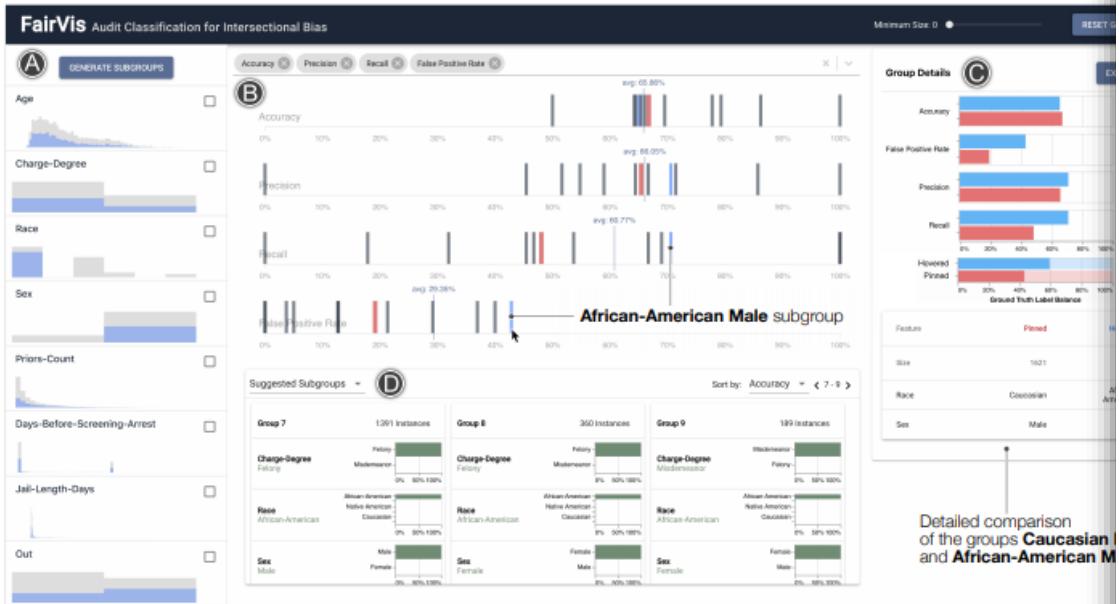


Figure 1: FAIRVIS integrates multiple coordinated views for discovering intersectional bias. Above, our user investigates intersectional subgroups of sex and race. **A.** The *Feature Distribution View* allows users to visualize each feature's distribution and generate subgroups. **B.** The *Subgroup Overview* lets users select various fairness metrics to see the global average performance and compare subgroups to one another, e.g., pinned **Caucasian Males** versus hovered **African-American Males**. The *Recall* and *False Positive Rate* show that for African-American Males, the model has relatively high recall but also the highest positive rate out of all subgroups of sex and race. **C.** The *Detailed Comparison View* lets users compare the details of two subgroups and investigate their class balances. Since the difference in False Positive Rates between Caucasian Males and African-American Males is far larger than their difference in base rates, a user suspects this part of the model merits further inquiry. **D.** The *Suggested Subgroups* and *Similar Subgroup View* shows suggested subgroups ranked by the worst performance in a given metric.

### FairSight: Visual Analytics for Fairness in Decision Making

Yongsu Ahn, Yu-Ru Lin

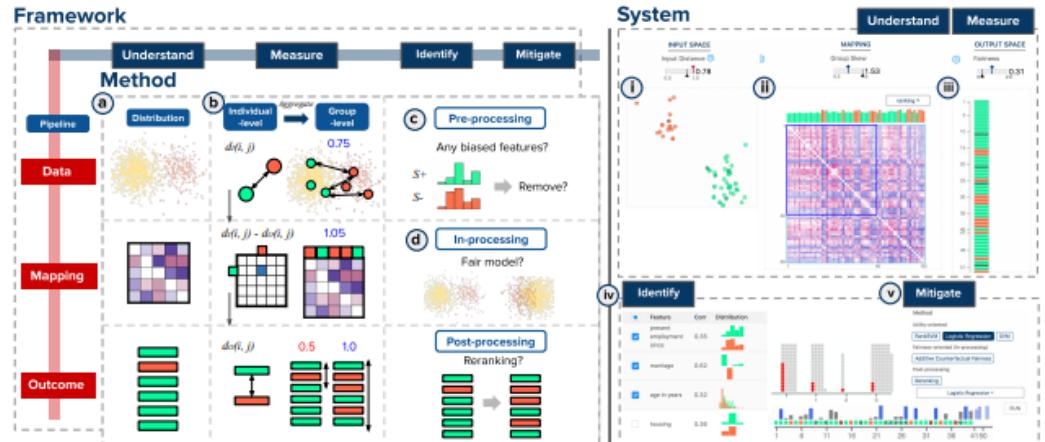


Fig. 1. We propose a design framework to protect individuals and groups from discrimination in algorithm-assisted decision making. A visual analytic system, *FairSight*, is implemented based on our proposed framework, to help data scientists and practitioners make fair decisions. The decision is made through ranking individuals who are either members of a protected group (orange bars) or a non-protected group (green bars). (a) The system provides a pipeline to help users understand the possible bias in a machine learning task as a *mapping* from the *input space* to the *output space*. (b) Different notions of fairness – *individual* fairness and *group* fairness – are measured and summarized numerically and visually. For example, the individual fairness is quantified by how pairwise distances between individuals are preserved through the mapping. The group fairness is quantified by the extent to which it leads to fair outcome distribution across groups, with (i) a 2D plot, (ii) a color-coded matrix, and (iii) a ranked-list plot capturing the pattern of potential biases. The system provides diagnostic modules to help (iv) identify and (v) mitigate biases through (c) investigating features before running a model, and (d) leveraging fairness-aware algorithms during and after the training step.

**Abstract**— Data-driven decision making related to individuals has become increasingly pervasive, but the issue concerning the potential discrimination has been raised by recent studies. In response, researchers have made efforts to propose and implement fairness measures and algorithms, but those efforts have not been translated to the real-world practice of data-driven decision making. As such, there is still an urgent need to create a viable tool to facilitate fair decision making. We propose *FairSight*, a visual analytic system to address this need; it is designed to achieve different notions of fairness in ranking decisions through the required actions – understanding, measuring, diagnosing and mitigating biases – that together lead to fairer decision making. Through a case study and user study, we demonstrate that the proposed visual analytic and diagnostic modules in the system are effective in understanding the fairness-aware decision pipeline and obtaining more fair outcomes.

**Index Terms**—Fairness in Machine Learning, Visual Analytic

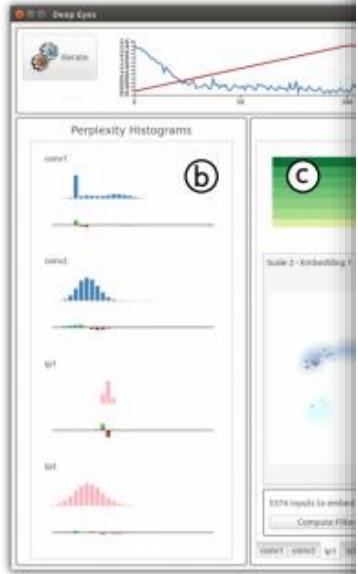
# Refinement

# Refinement

## Single Model

## DeepEyes: Progressive Visual Analytics for Designing Deep Neural Networks

Nicola Pezzotti, Thomas Höllt, Jar



**Fig. 1. DeepEyes** is a Progressive Visual on the training is given by the commonly that allows the detection of stable layers. filters are detected in the Activation Heatmap relationships among the filters in a layer a

**Abstract**—Deep neural networks are now classifiers, where features are handcrafted. Handcrafting the features, it is now the net the number of layers or the number of filter design guidelines exist, designing a neural network due to the large datasets used for training the design of neural networks during training, a stable set of patterns and, therefore, are of superfluous filters or layers, and information through multiple use cases, showing how

**Index Terms**—Progressive visual analytic

# RetainVis: Visual Analytics with Interpretable and Interactive Recurrent Neural Networks on Electronic Medical Records

Bum Chul Kw  
Young Bin Kim

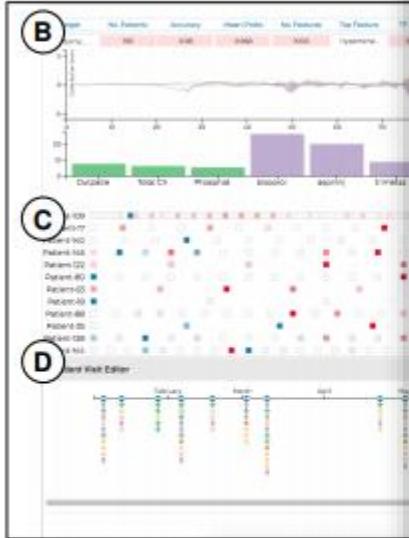


Fig. 1. A screenshot of RetainVis consists of a summary view (right) of patients. (B) Patient list view (left) showing a list of individual patients in a row of rectangles. (C) Patient detail view. Users can open (D) *Patient Editor* to conduct further analysis.

**Abstract**— We have recently seen many systems (EMRs), which contain histories of patients' states of patients. Despite the strong performance of these systems, they lack the ability for particular prediction. Such *black-box* nature of these systems makes it difficult to establish methods to interactively leverage their predictions. In this paper, we present our design study aims to provide a visual interface for medical experts, artificial intelligence scientists, and domain experts, we design, implement, and evaluate a system that can interactively leverage the predictions of an interactive RNN-based model called RNN-EMR. Our study shows the effective use of the system for particular prediction of the patients' states using EMRs of patients with heart diseases.

# Towards Better Analysis of Deep Convolutional Neural Networks

Mengchen Liu, Jiaxin Shi, Zhen Li, Chongxuan Li, Jun Zhu, Shixia Liu

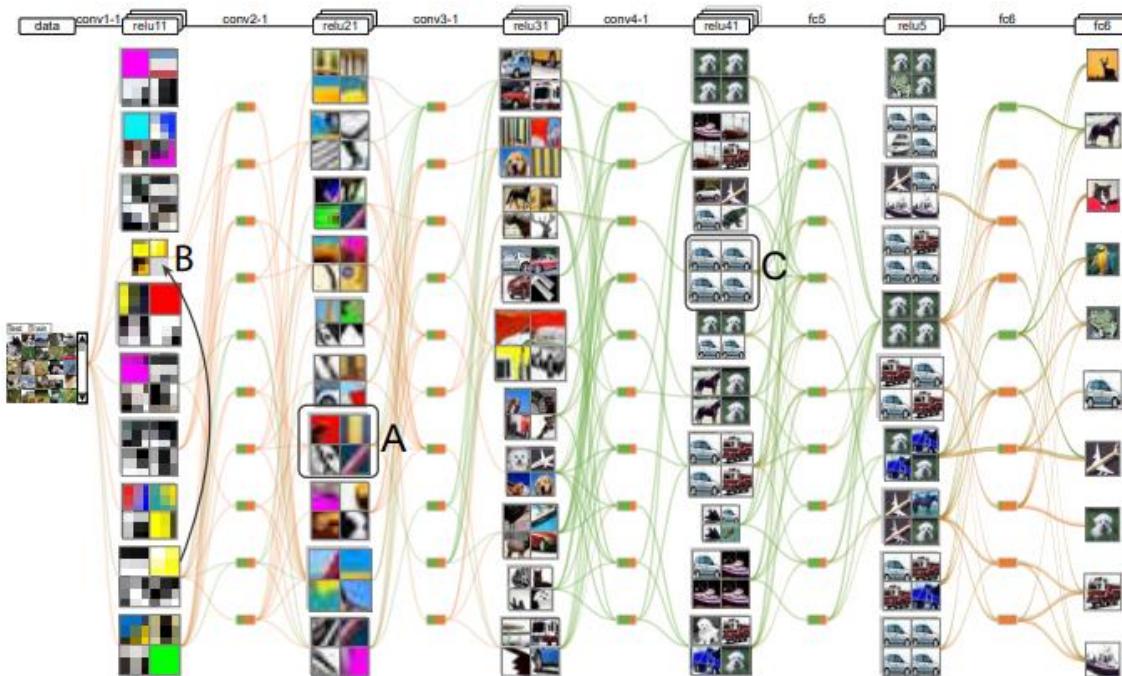


Fig. 1. CNNVis, a visual analytics toolkit that helps experts understand, diagnose, and refine deep CNNs

**Abstract**— Deep convolutional neural networks (CNNs) have achieved breakthrough performance in many pattern recognition tasks such as image classification. However, the development of high-quality deep models typically relies on a substantial amount of

# DeepCompare: Visual and Interactive Comparison of Deep Learning Model Performance

Sugeerth Murugesan  
University of California, Davis

Sana Malik, Fan Du, and Eunyee Koh  
Adobe Research

**Abstract**—Deep learning models have become the state-of-the-art in text sentiment analysis to facial image recognition. However, certain models perform better than others or how one model outperforms another is often difficult yet critical for increasing their effectiveness, prediction accuracy, and enabling fairness. Traditional metrics of model efficacy, such as accuracy, precision, and recall provide a quantitative view of model performance; however, the qualitative intricacies of why one model outperforms another are hidden. In this paper, we interview machine learning experts to understand their evaluation and comparison workflow. From this, we propose a visual analytic approach, DeepCompare, to systematically compare two deep learning models, in order to provide insight into the model's strengths and weaknesses, and interactively assess tradeoffs between two such models. The system allows users to evaluate model results, identify and compare activation patterns, and

Tuan Manh La  
Purdue University

## Progressive Learning of Topic Modeling Parameters: A Visual Analytics Framework

Mennatallah El-Assady<sup>1,2</sup>, Rita Sevastjanova<sup>1</sup>, Fabian Sperrle<sup>1</sup>, Daniel Keim<sup>1</sup>, and Christopher Collins<sup>2</sup>

<sup>1</sup>University of Konstanz, Germany

<sup>2</sup>University of Ontario Institute of Technology, Canada

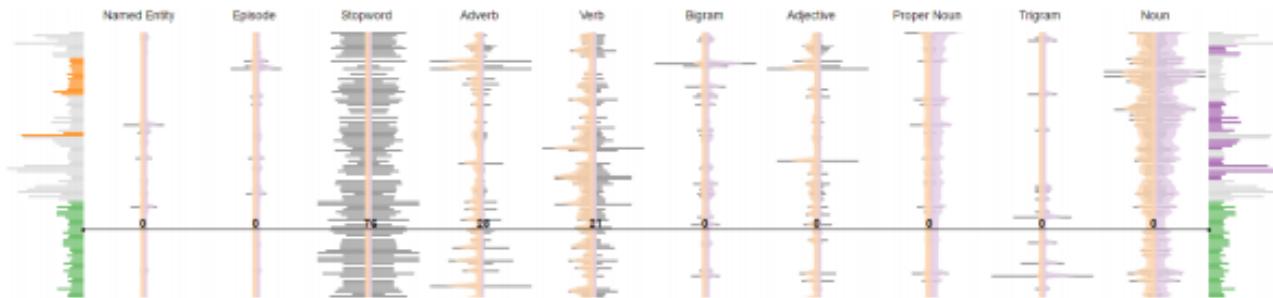


Fig. 1. *Parameter Distribution View* using comparative bar charts. This compact visualization technique enhances the comparison of two parameter distributions using mirrored bar-charts as a baseline and two asymmetrical violin-style plots as distribution estimates. The plots are scaled using the ratio between the two compared assortments (on both sides). The larger value is scaled to the full width of the baseline and the smaller value is scaled proportionally. This figure depicts the comparison of the utterance descriptor features of the second US presidential debate between Obama and Romney in 2012. All utterances are sorted according to their topic coherence.

**Abstract**— Topic modeling algorithms are widely used to analyze the thematic composition of text corpora but remain difficult to interpret and adjust. Addressing these limitations, we present a modular visual analytics framework, tackling the understandability and adaptability of topic models through a user-driven reinforcement learning process which does not require a deep understanding of the underlying topic modeling algorithms. Given a document corpus, our approach initializes two algorithm configurations based on a parameter space analysis that enhances document separability. We abstract the model complexity in an interactive visual workspace for exploring the automatic matching results of two models, investigating topic summaries, analyzing parameter distributions, and reviewing documents. The main contribution of our work is an iterative decision-making technique in which users provide a document-based relevance feedback that allows the framework to converge to a user-endorsed topic distribution. We also report feedback from a

# Human-Centered Tools for Coping with Imperfect Algorithms During Medical Decision-Making

Carrie J. Cai, Emily Reif, Narayan Hegde, Jason Hipp, Been Kim, Daniel Smilkov,  
 Martin Wattenberg, Fernanda Viegas, Greg S. Corrado, Martin C. Stumpe, Michael Terry  
 Google Brain, Google Health  
 Mountain View, CA  
 {cjcai,ereif,hegde,hipp,beenkim,smilkov,wattenberg,viegas,ccorrado,mstumpe,m

## ABSTRACT

Machine learning (ML) is increasingly being used in image retrieval systems for medical decision making. One application of ML is to retrieve visually similar medical images from past patients (e.g. tissue from biopsies) to reference when making a medical decision with a new patient. However, no algorithm can perfectly capture an expert's ideal notion of similarity for every case: an image that is algorithmically determined to be similar may not be medically relevant to a doctor's specific diagnostic needs. In this paper, we identified the needs of pathologists when searching for similar images retrieved using a deep learning algorithm, and developed tools that empower users to cope with the search algorithm on-the-fly, communicating what types of similarity are most important at different moments in time. In two evaluations with pathologists, we found that these refinement tools increased the diagnostic utility of images found and increased user trust in the algorithm. The tools were preferred over a traditional interface, without a loss in diagnostic accuracy. We also observed that users adopted new strategies when using refinement tools, re-purposing them to test and understand the underlying algorithm and to disambiguate ML errors from their own errors. Taken together, these findings inform future human-ML collaborative systems for expert decision-making.

## CCS CONCEPTS

- Human-centered computing → Human computer interaction (HCI);

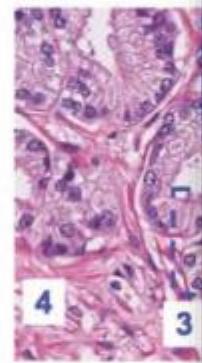


Figure 1: Medical image features, such as cellular interaction between core (4), and many more. It can search algorithm to perform similarity, because what's important differs from dependent.

## ACM Reference Format:

Carrie J. Cai, Emily Reif, Narayan Hegde, Jason Hipp, Been Kim, Daniel Smilkov, Martin Wattenberg, Fernanda Viegas, Greg S. Corrado, Martin C. Stumpe, Michael Terry. 2019. Tools for Coping with Imperfect Algorithms During Medical Decision-Making. In *CHI Conference on Human Factors in Computing Systems Proceedings (CHI 2019)*. ACM, New York, NY, USA, 1–12. 3300234

## Semantic Concept Spaces: Guided Topic Model Refinement using Word-Embedding Projections

Mennatallah El-Assady<sup>1,2</sup>, Rebecca Kehlbeck<sup>1</sup>, Christopher Collins<sup>2</sup>, Daniel Keim<sup>1</sup>, and Oliver Deussen<sup>1</sup>

<sup>1</sup> University of Konstanz, Germany.

<sup>2</sup> University of Ontario Institute of Technology, Canada.



Fig. 1: Guided relevance feedback for the targeted refinement of incoherent areas in the *Semantic Concept Space*. This user guidance component tours through the space and highlights potentially uncertain areas, suggesting a recommended action for refinement.

**Abstract**— We present a framework that allows users to incorporate the semantics of their domain knowledge for topic model refinement while remaining model-agnostic. Our approach enables users to (1) *understand* the semantic space of the model, (2) *identify* regions of potential conflicts and problems, and (3) *readjust* the semantic relation of concepts based on their understanding, directly influencing the topic modeling. These tasks are supported by an interactive visual analytics workspace that uses word-embedding projections to define *concept regions* which can then be refined. The user-refined concepts are independent of a particular document collection and can be transferred to related corpora. All user interactions within the concept space directly affect the semantic relations of the underlying vector space model which, in turn, change the topic modeling. In addition to direct manipulation, our system guides the users' decision-making process through recommended interactions that point out potential improvements. This targeted refinement aims at minimizing the feedback required for an efficient human-in-the-loop process. We confirm the improvements achieved through our approach in two

# Beyond VIS

## Minions, Sheep, and Fruits: Metaphorical Narratives to Explain Artificial Intelligence and Build Trust

Wolfgang Jentner<sup>1</sup>

### *Shall we play? – Extending the Visual Analytics Design Space through Gameful Design Concepts*

Rita Sevestjanova<sup>\*</sup>

University of Konstanz



Fig. 1. A transitive diagram showing the complexity of game design.

**Abstract**— Advances in game design are bringing along the expectation that the underlying artificial intelligence should be effective and helpful. Therefore, we propose to involve domain experts and domain experts to build complex models to explain the building and an adequate model is known to provide trust in fields and theories.

Figure 1: We describe a model for **When** does the challenging task occur, the users do to design an engaging game. **Why** do people do those challenging tasks have impact, and to be accepted. **What** are the game mechanics such

### Going beyond Visualization: Verbalization as Complementary Medium to Explain Machine Learning Models

Rita Sevestjanova<sup>1</sup>, Fabian Beck<sup>2</sup>, Basil Ell<sup>3</sup>, Cagatay Turkay<sup>4</sup>, Rafael Henkin<sup>4</sup>, Miriam Butt<sup>1</sup>, Daniel Keim<sup>1</sup>, Mennatallah El-Assady<sup>1,5</sup>

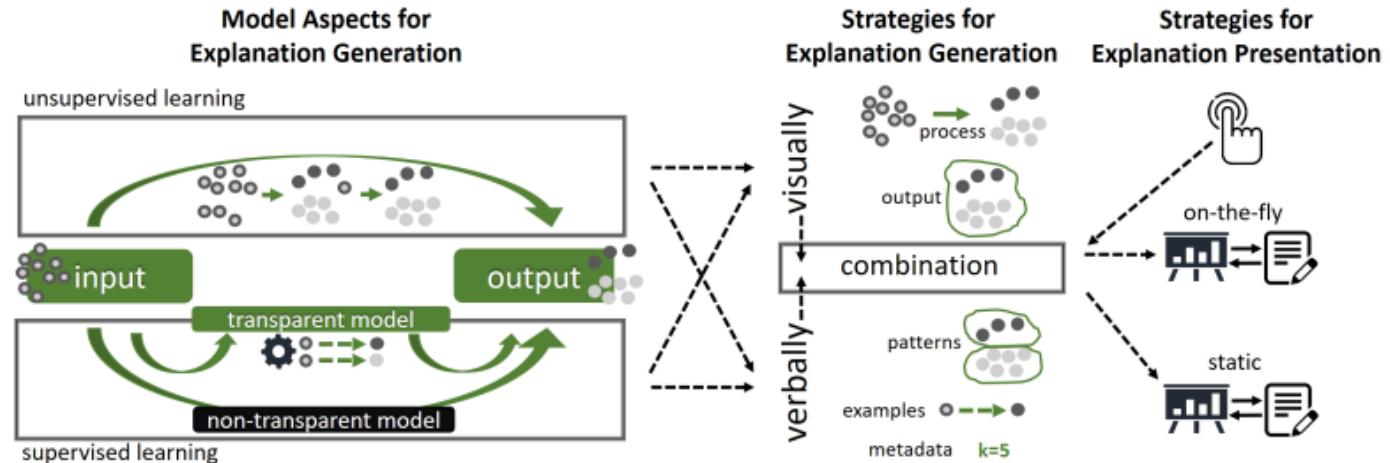
<sup>1</sup>University of Konstanz, Germany

<sup>2</sup>University of Duisburg-Essen, Germany

<sup>3</sup>Bielefeld University, Germany

<sup>4</sup>City University of London, UK

<sup>5</sup>University of Ontario Institute of Technology, Canada





The framework for explainable AI and interactive machine learning. Making XAI accessible.

<https://explainer.ai/>

## explAIner: A Visual Analytics Framework for Interactive and Explainable Machine Learning

Thilo Spinner, Udo Schlegel, Hanna Schäfer, and Mennatallah El-Assady

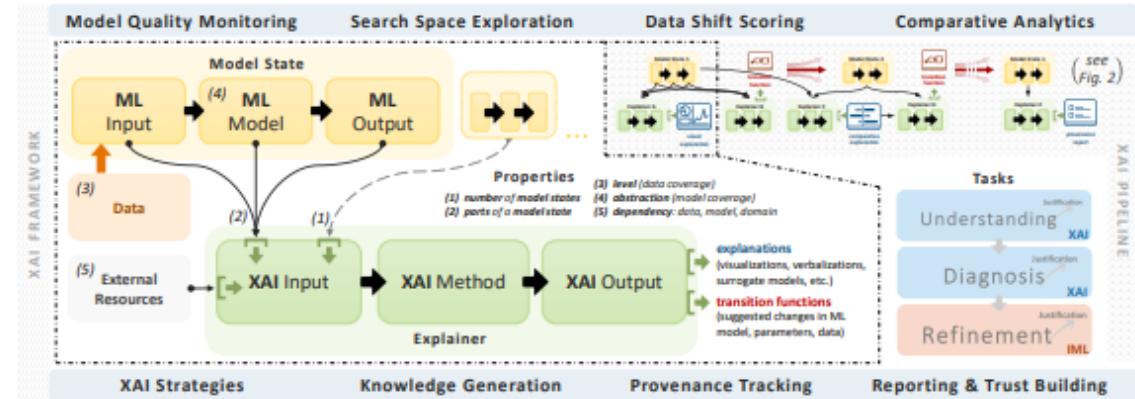


Fig. 1: Close-up view of an *explainer*, the main building-block used to construct an iterative *XAI pipeline* for the *understanding*, *diagnosis*, and *refinement* of ML models. Explainers have five properties; they take one or more model states as input, applying an XAI method, to output an explanation or a transition function. Global monitoring and steering mechanisms expand the pipeline to the full *XAI framework*, supporting the overall workflow by guiding, steering, or tracking the explainers during all steps.

**Abstract**— We propose a framework for interactive and explainable machine learning that enables users to (1) understand machine learning models; (2) diagnose model limitations using different explainable AI methods; as well as (3) refine and optimize the models. Our framework combines an iterative XAI pipeline with eight global monitoring and steering mechanisms, including quality monitoring, provenance tracking, model comparison, and trust building. To operationalize the framework, we present explAIner, a visual analytics system for interactive and explainable machine learning that instantiates all phases of the suggested pipeline within the commonly used TensorBoard environment. We performed a user-study with nine participants across different expertise levels to examine their perception of our workflow and to collect suggestions to fill the gap between our system and framework. The evaluation confirms that our tightly integrated system leads to an informed machine learning process while disclosing opportunities for further extensions.

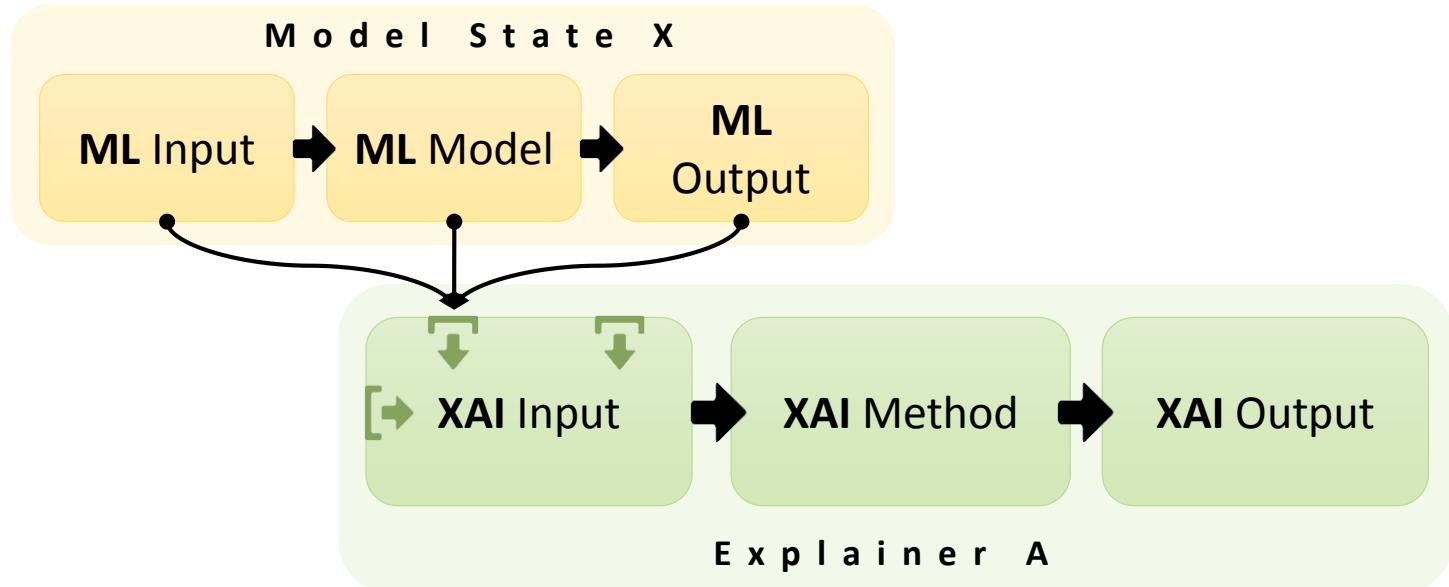
**Index Terms**— Explainable AI, Interactive Machine Learning, Deep Learning, Visual Analytics, Interpretability, Explainability

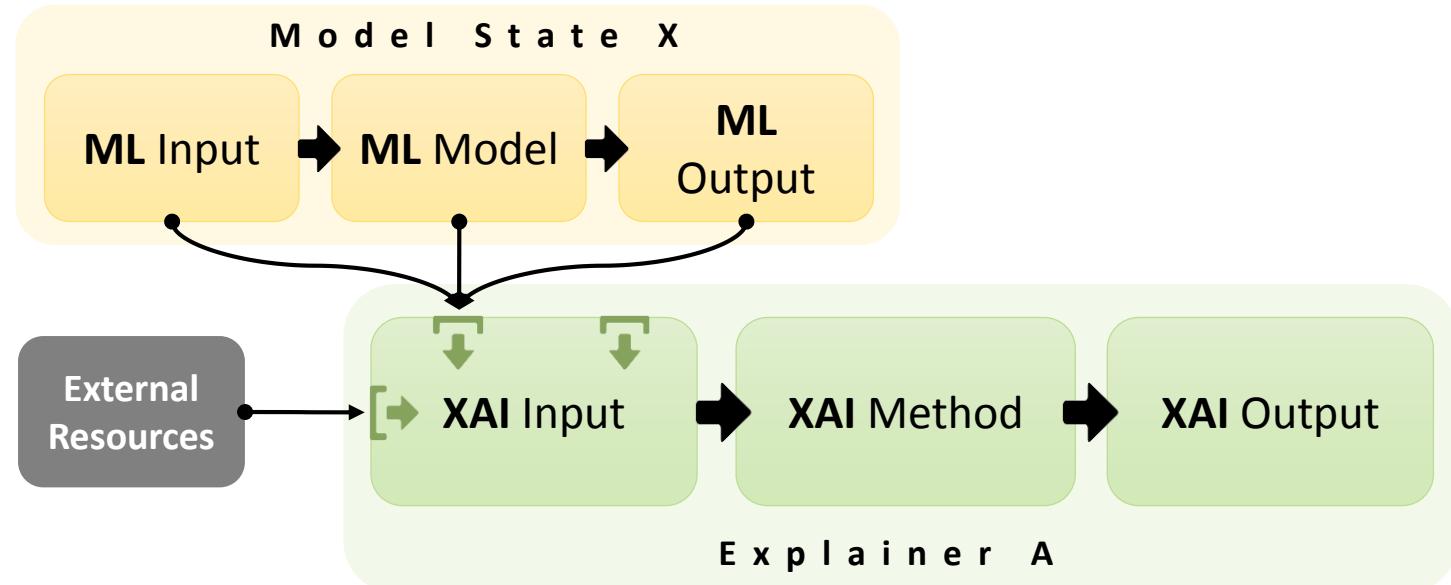


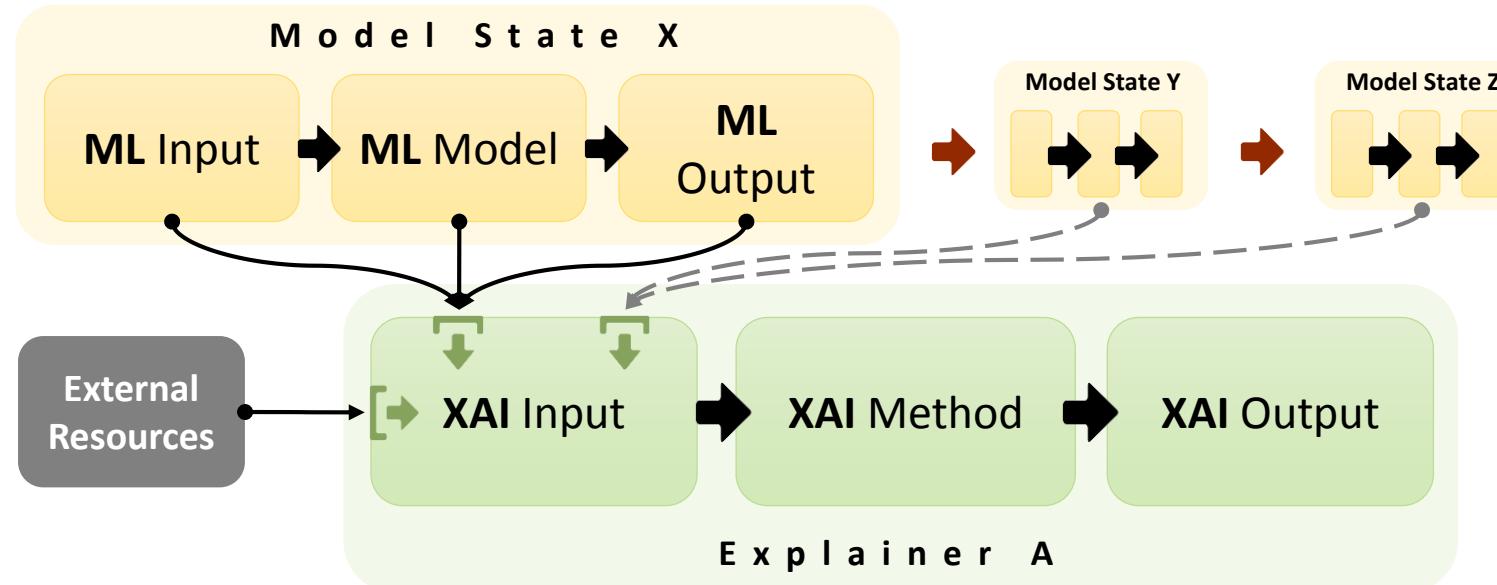
# XAI as a Process

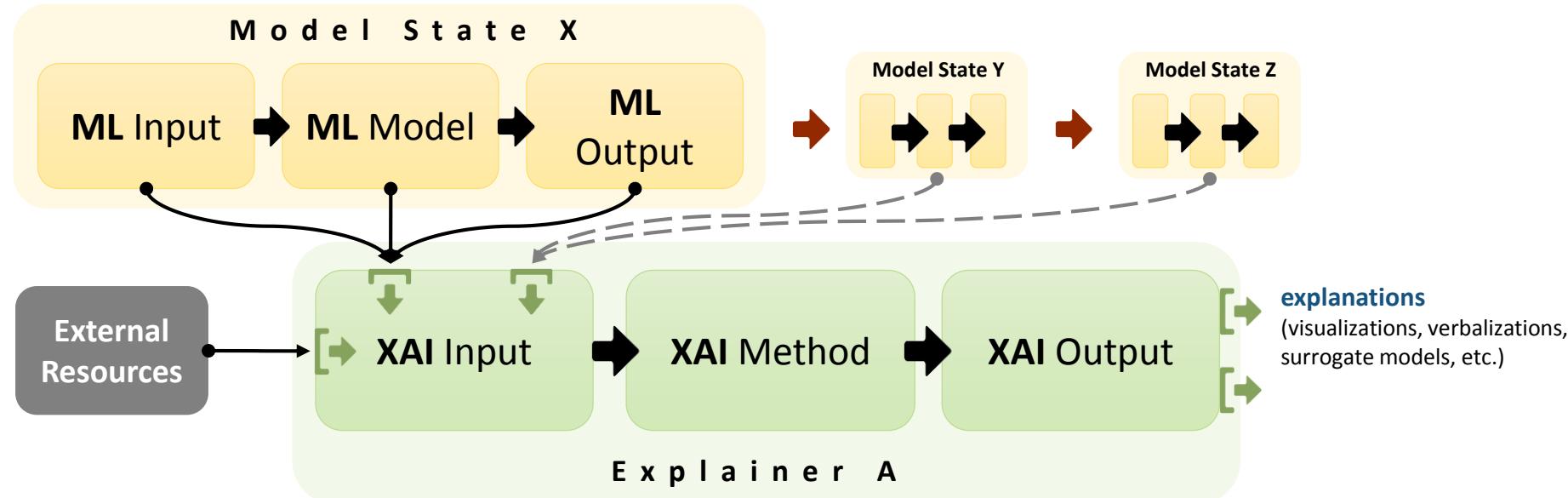
M o d e l   S t a t e   X

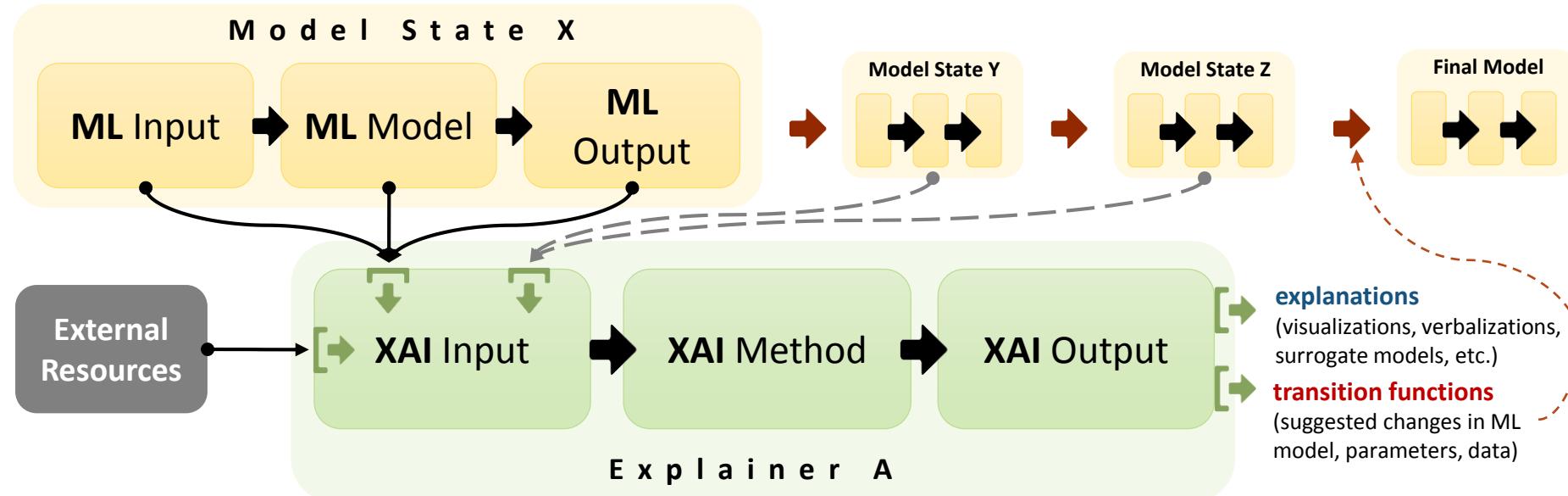


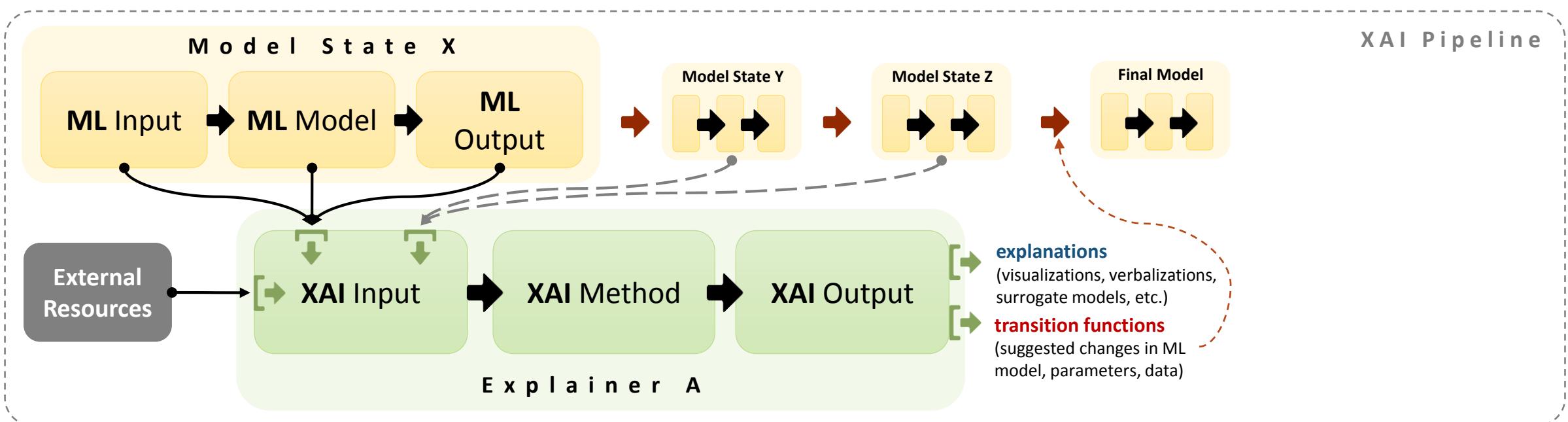


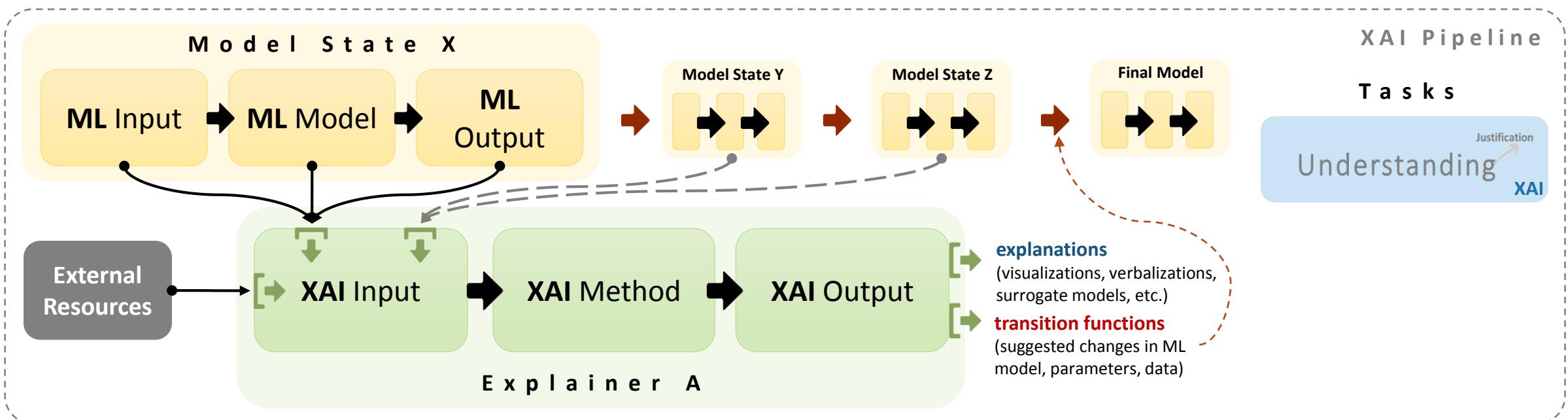


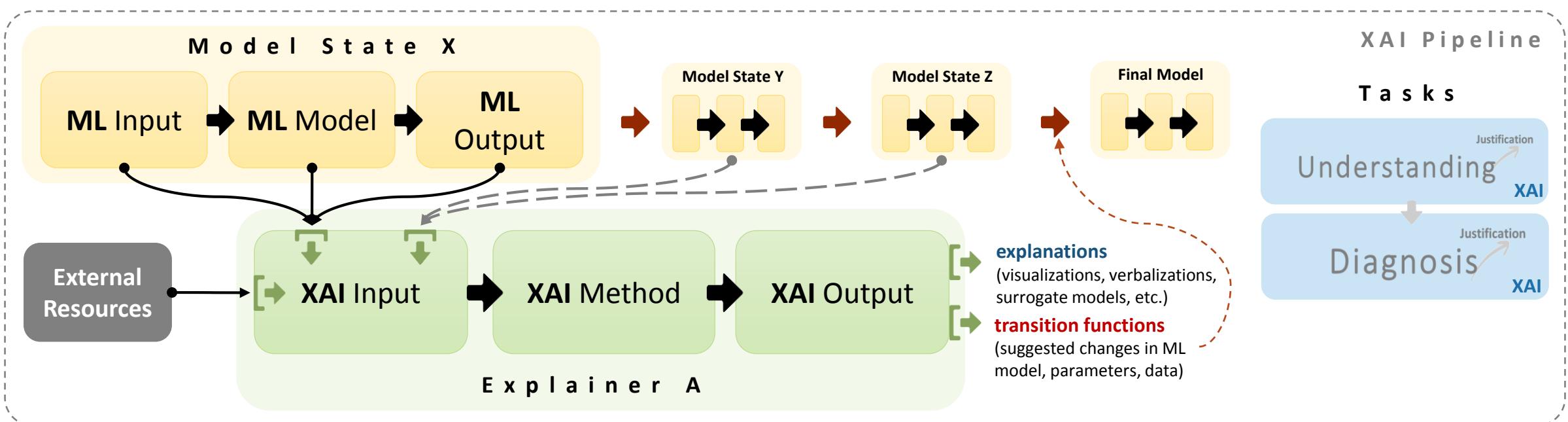


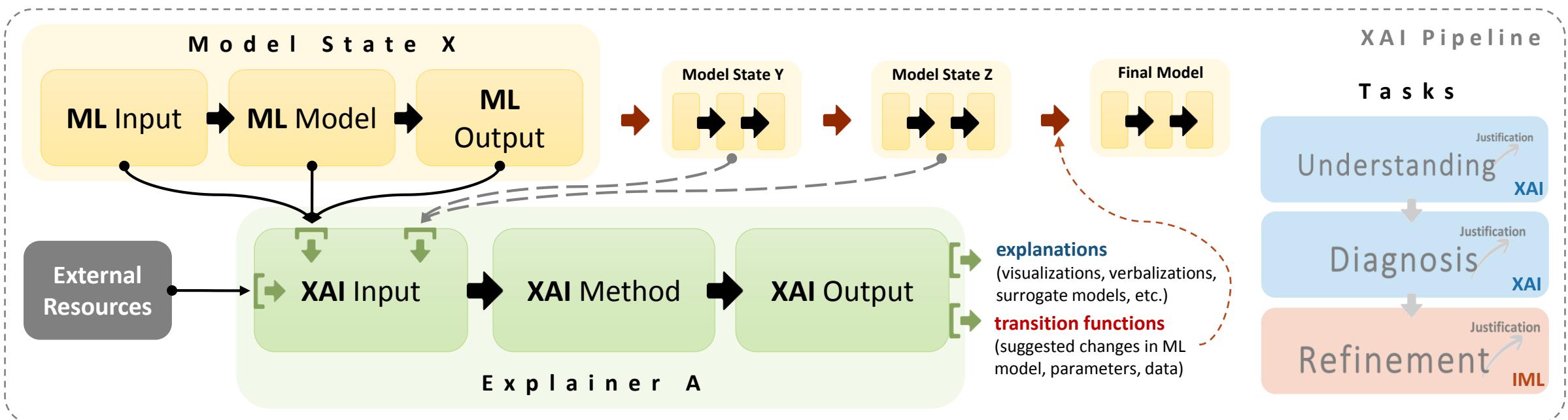




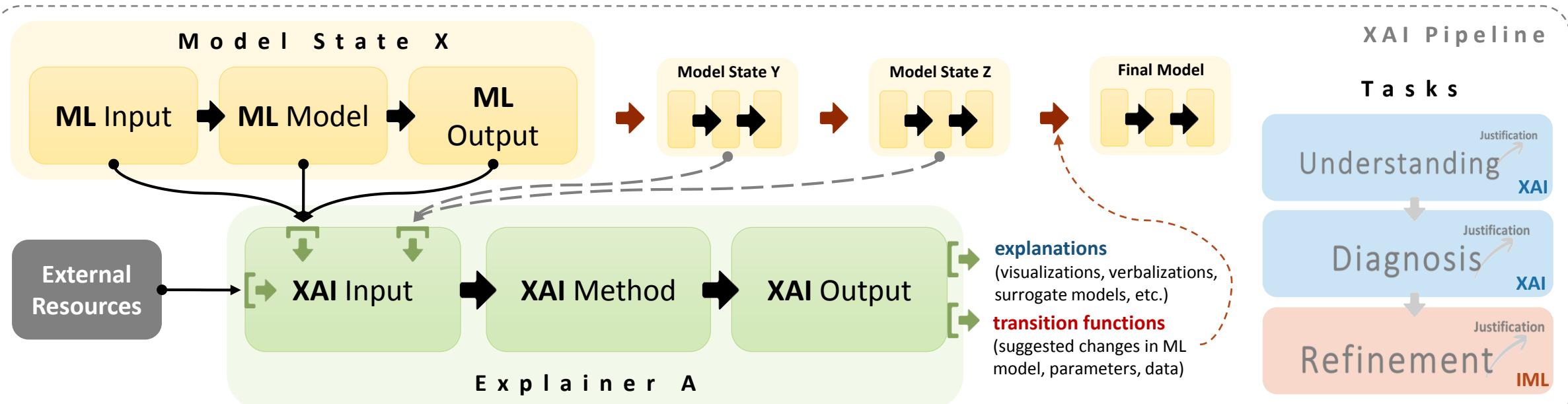


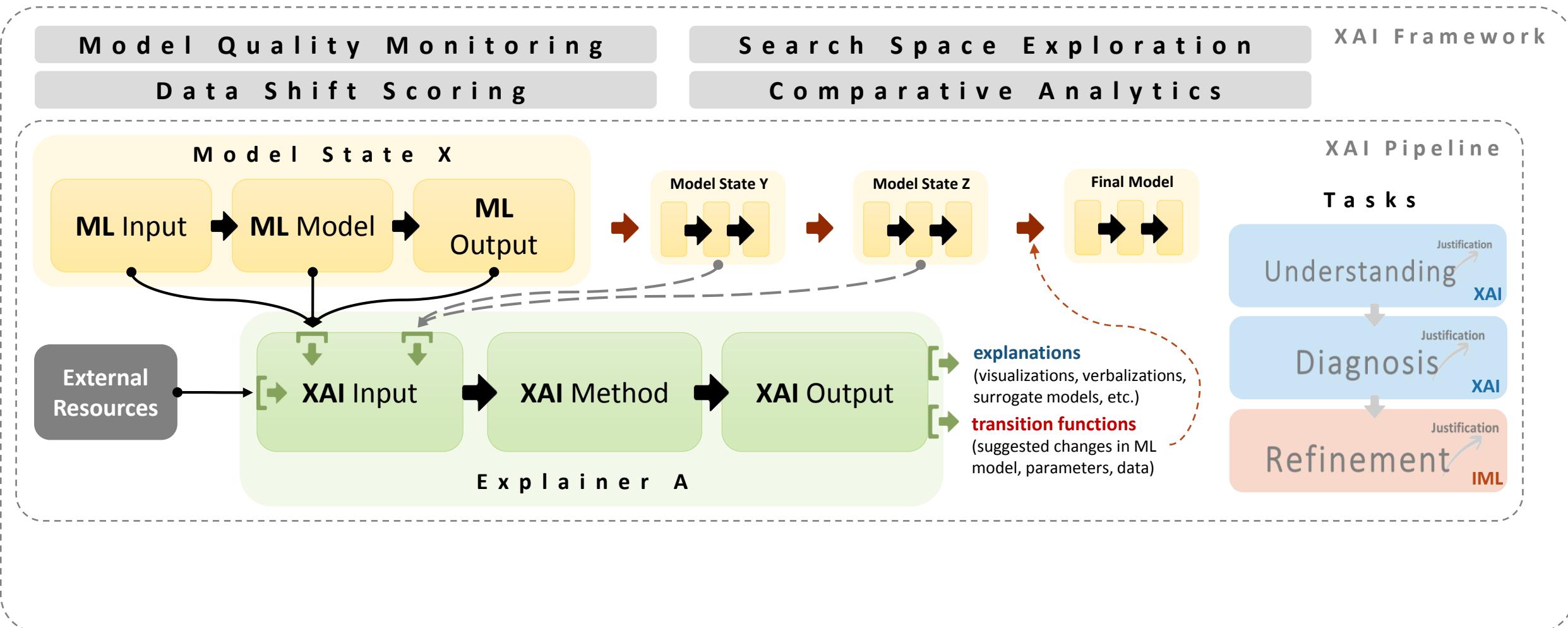


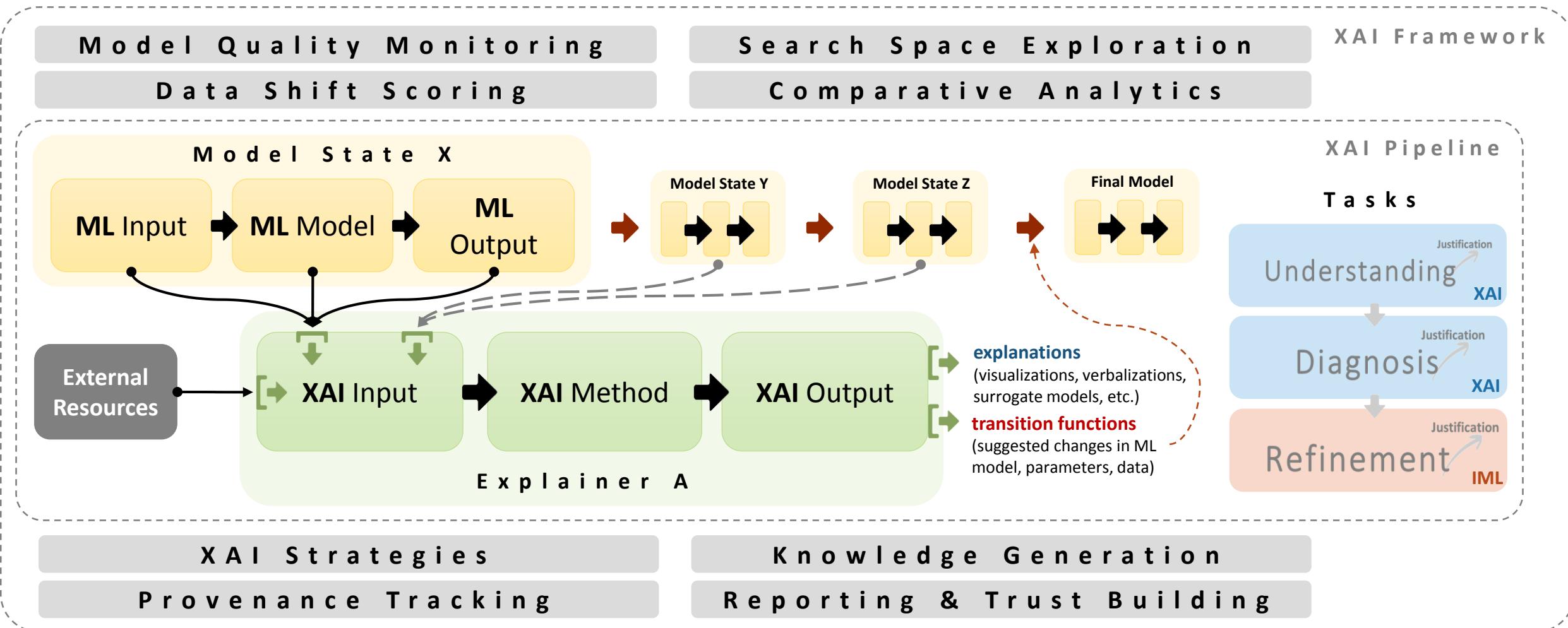




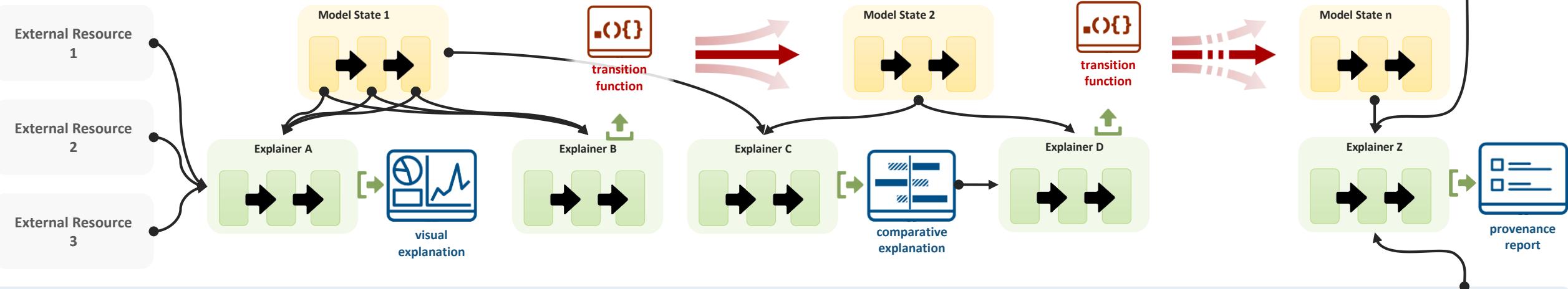
## XAI Framework



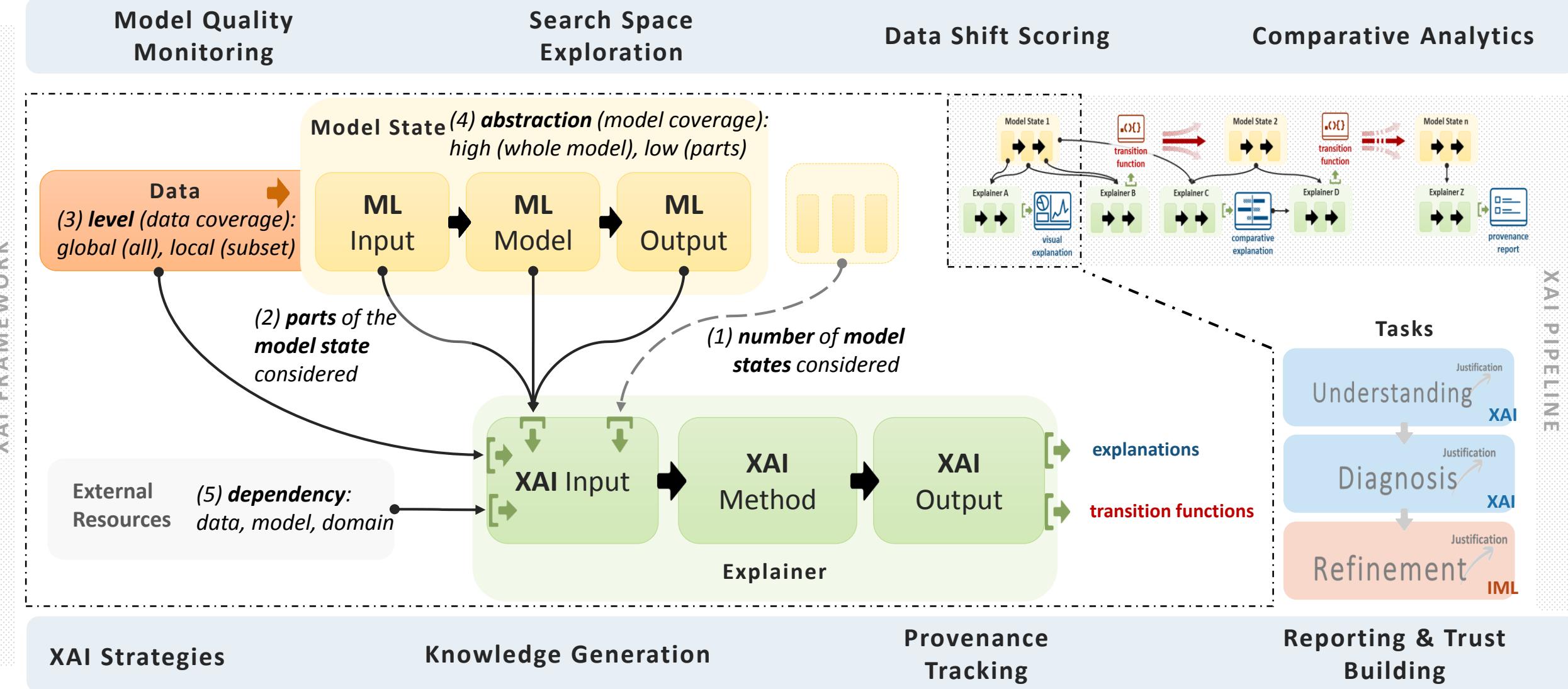




## Model Quality Monitoring

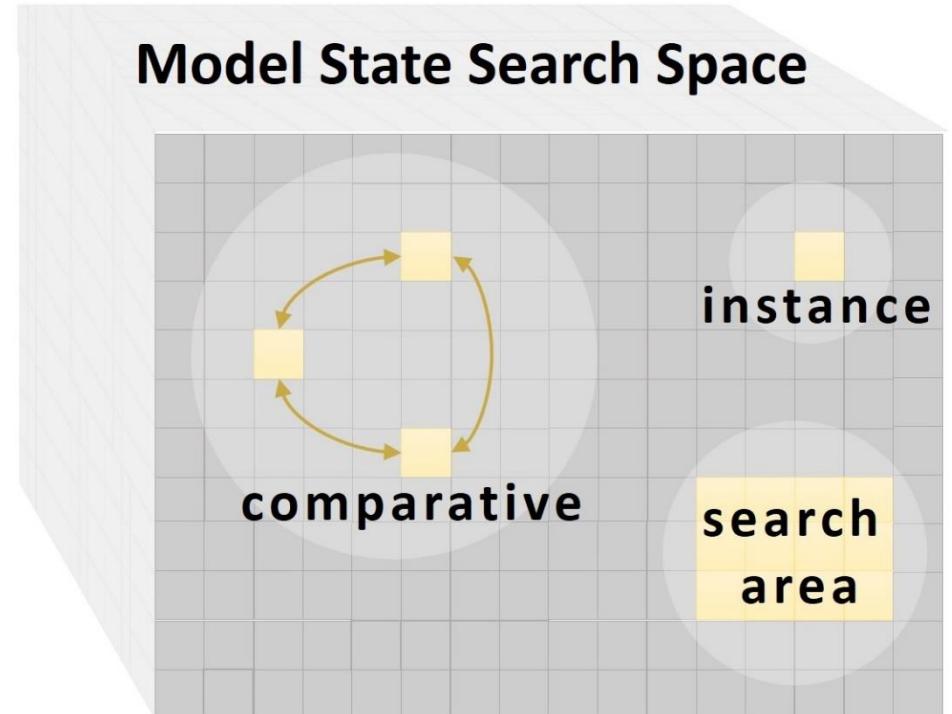


## Provenance Tracking



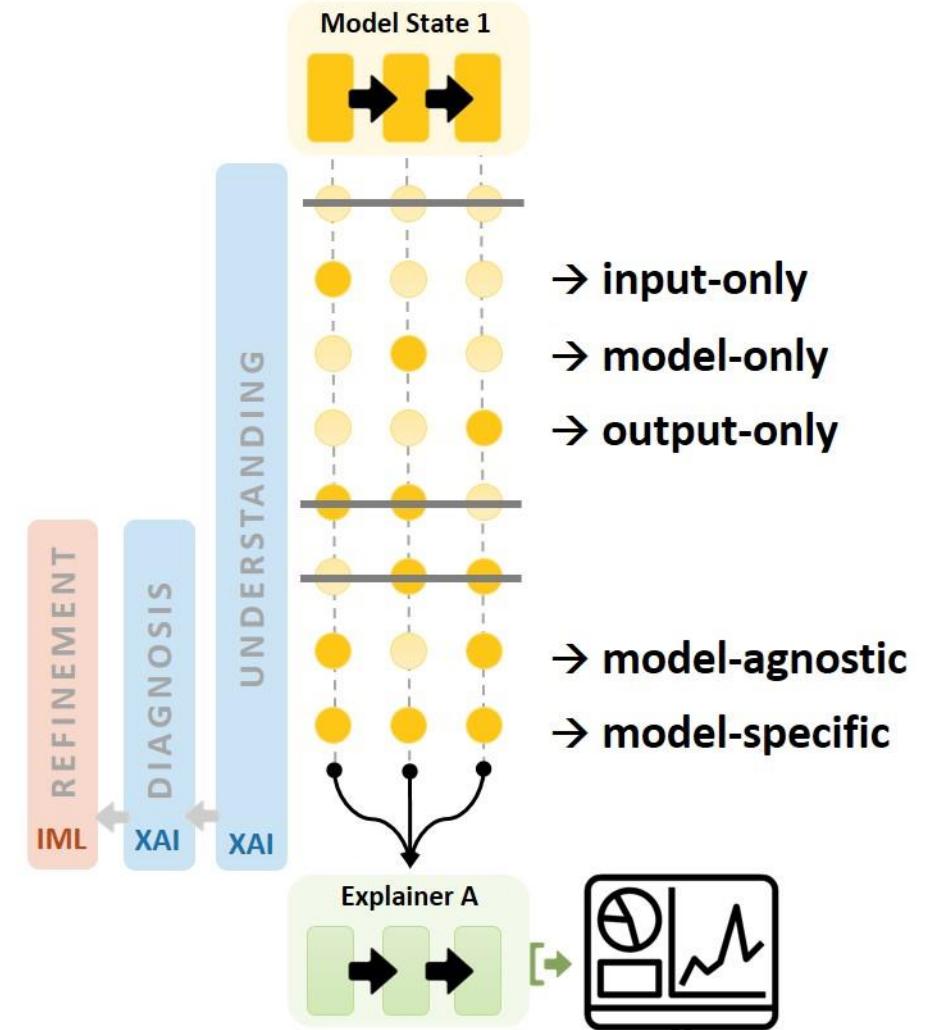
# Explainer Properties

- (1) number of model states considered**
- (2) parts of the model state considered**
- (3) level (data coverage): global, local**
- (4) abstraction (model coverage): high, low**
- (5) dependency: data, model, domain**



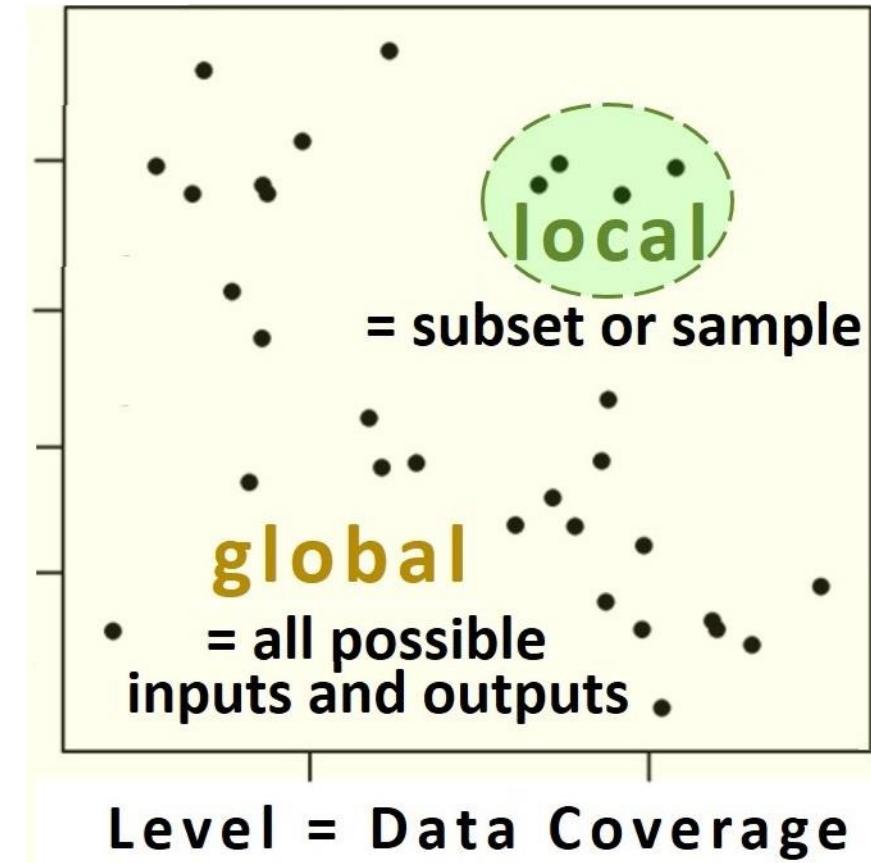
# Explainer Properties

- (1) number of model states considered**
- (2) parts of the *model state* considered**
- (3) level (data coverage): global, local**
- (4) abstraction (model coverage): high, low**
- (5) dependency: data, model, domain**



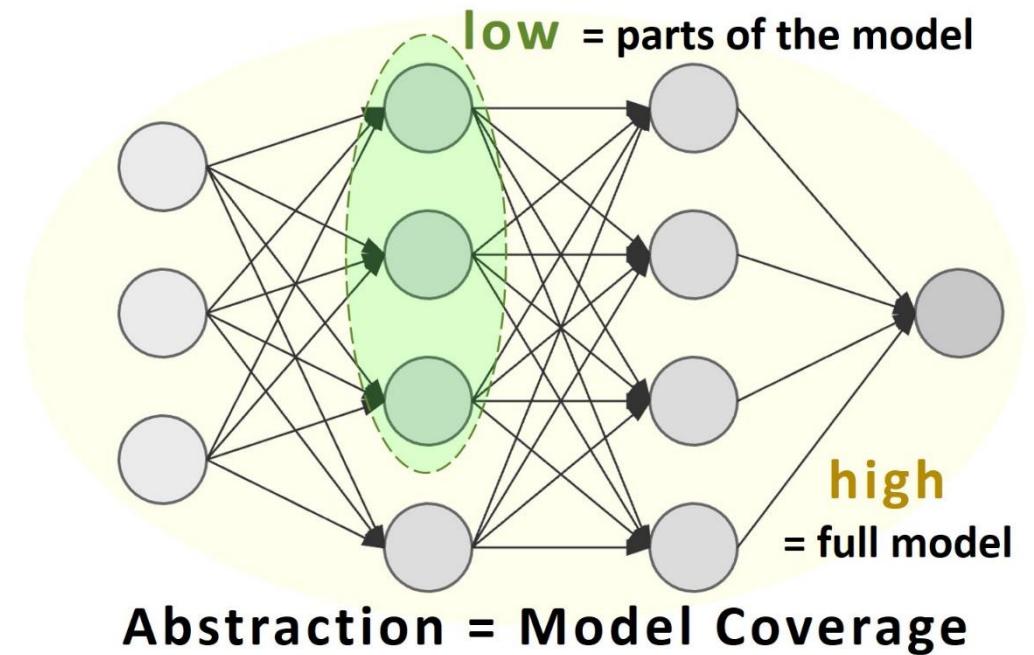
# Explainer Properties

- (1) *number of model states considered*
- (2) *parts of the model state considered*
- (3) **level** (*data coverage*): *global*, *local*
- (4) *abstraction* (*model coverage*): *high*, *low*
- (5) *dependency*: *data*, *model*, *domain*



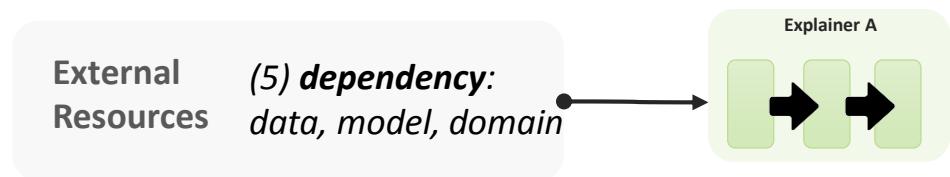
# Explainer Properties

- (1) number of model states considered*
- (2) parts of the model state considered*
- (3) level (data coverage): global, local*
- (4) abstraction (model coverage): high, low*
- (5) dependency: data, model, domain*



# Explainer Properties

- (1) number of model states considered**
- (2) parts of the model state considered**
- (3) level (data coverage): global, local**
- (4) abstraction (model coverage): high, low**
- (5) dependency: data, model, domain**



## Global Surrogate Models

-  Lime
-  Anchors
-  Cav

## Global Gradient Propagation Methods

-  Grad\*Input
-  IntGrad
-  Gradient
-  SmoothGrad
-  Input Times Gradient
-  Integrated Gradients
-  DeepLift
-  Saliency
-  DeConvNet

## Global Layer-wise Relevance Propagation

-  z-LRP
-  e-LRP (fast)
-  e-LRP (slow)

## Model

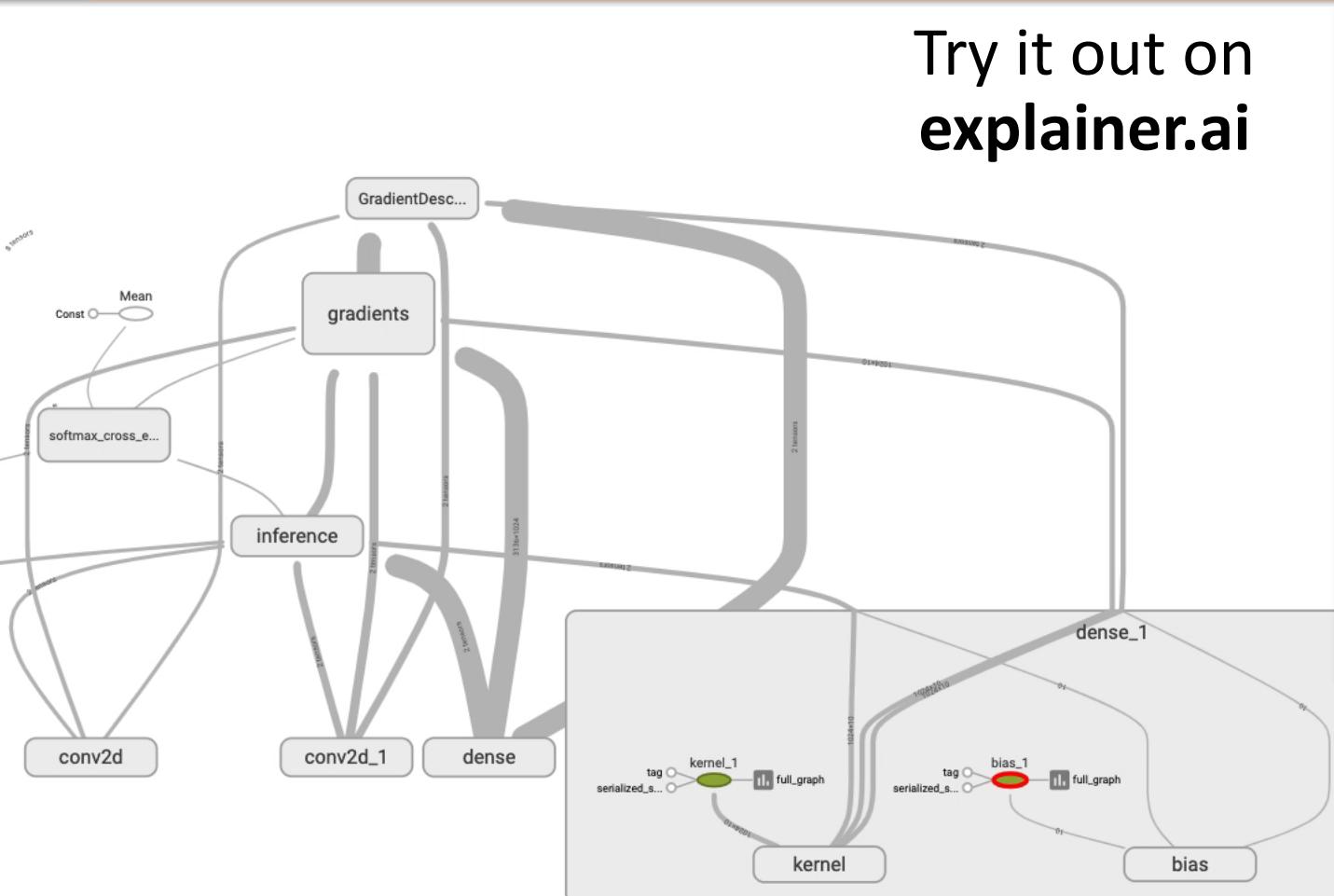
-  FunctionExplainer
-  LossDecay

## Graph Flow

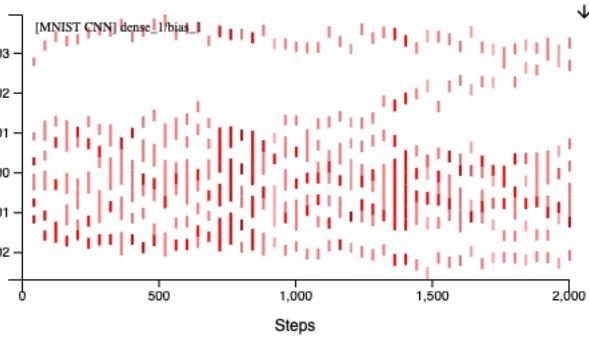
-  ImageInterpreter
-  InformationFlow

## Weights

-  HistoTrend
-  MinMax
-  DeadWeight
-  SaturatedWeight

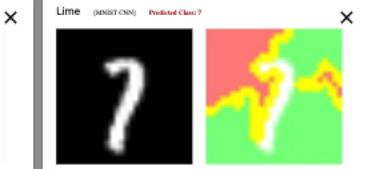
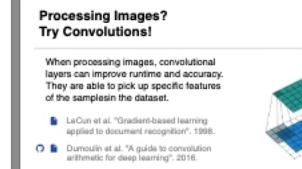
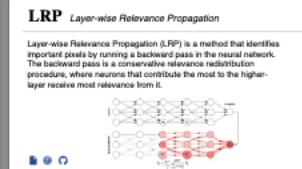
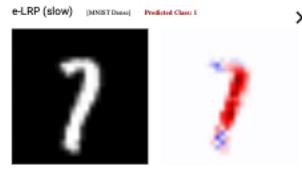
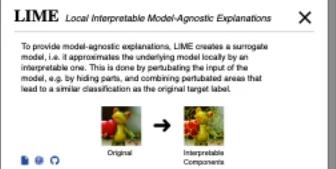
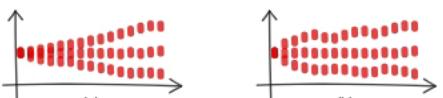


# Try it out on [explainer.ai](https://explainer.ai)



## HistoTrend Histograms, over Time

The HistoTrend-Tool plots the Histograms of a tensor, one for each logged step. It shows how the distribution of values changes over time and can reveal interesting patterns, such as binning (a) or weight oscillations (b).

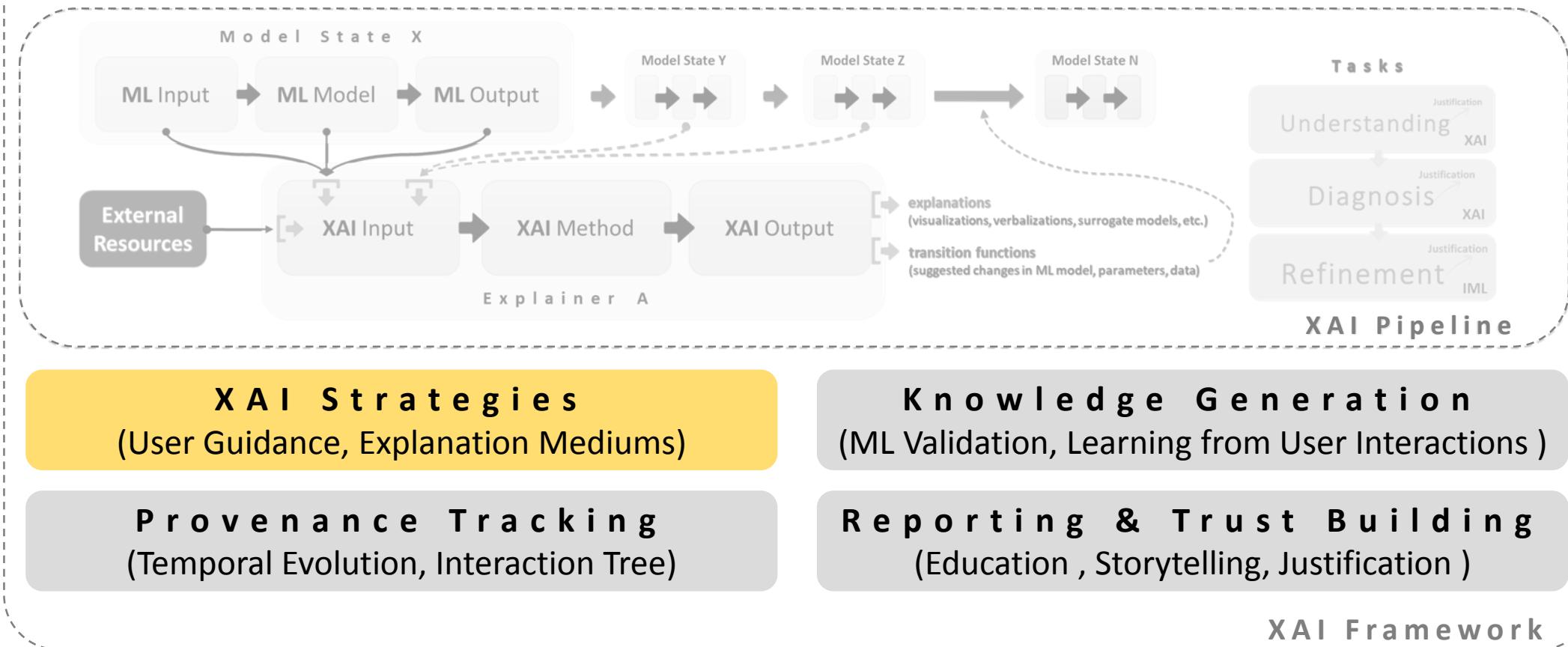


**Model Quality Monitoring**  
(Metrics, Bias, Uncertainty, Performance)

**Data Shift Scoring**  
(Training -- Testing Data Continuum )

**Search Space Exploration**  
(Speculative Execution, Targeted Optimization)

**Comparative Analytics**  
(Model Selection, Recommendation )



Global Monitoring and Steering Mechanisms



Pedagogy



Storytelling



Dialog & Argumentation

# Towards XAI

## Structuring the Processes of Explanations



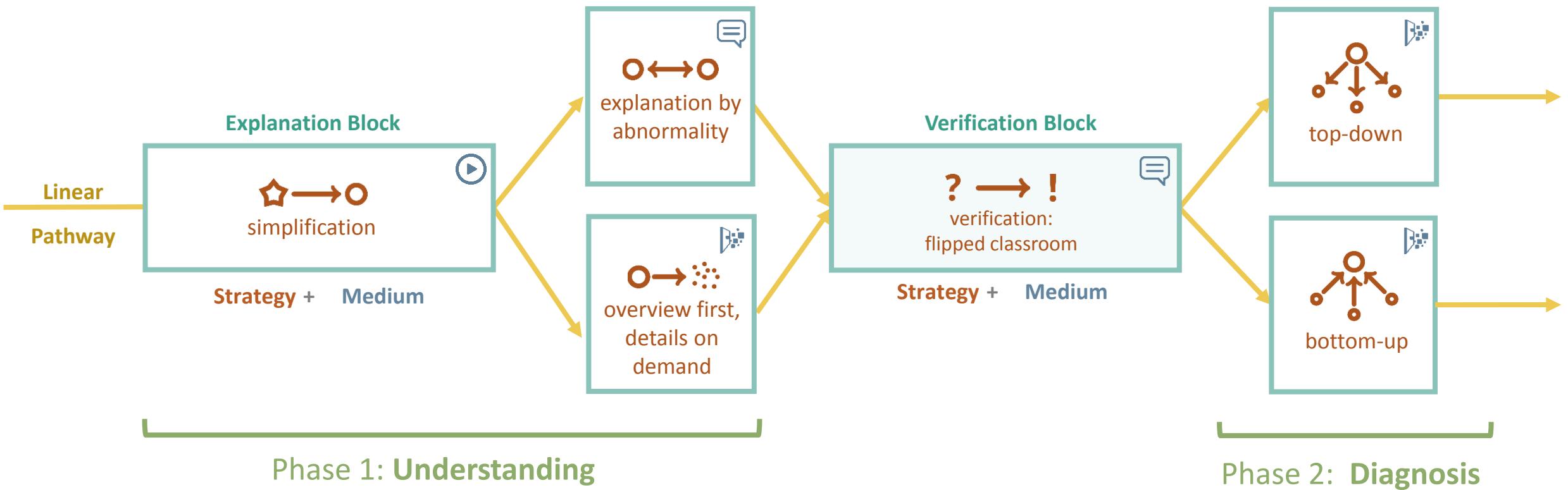
Programming



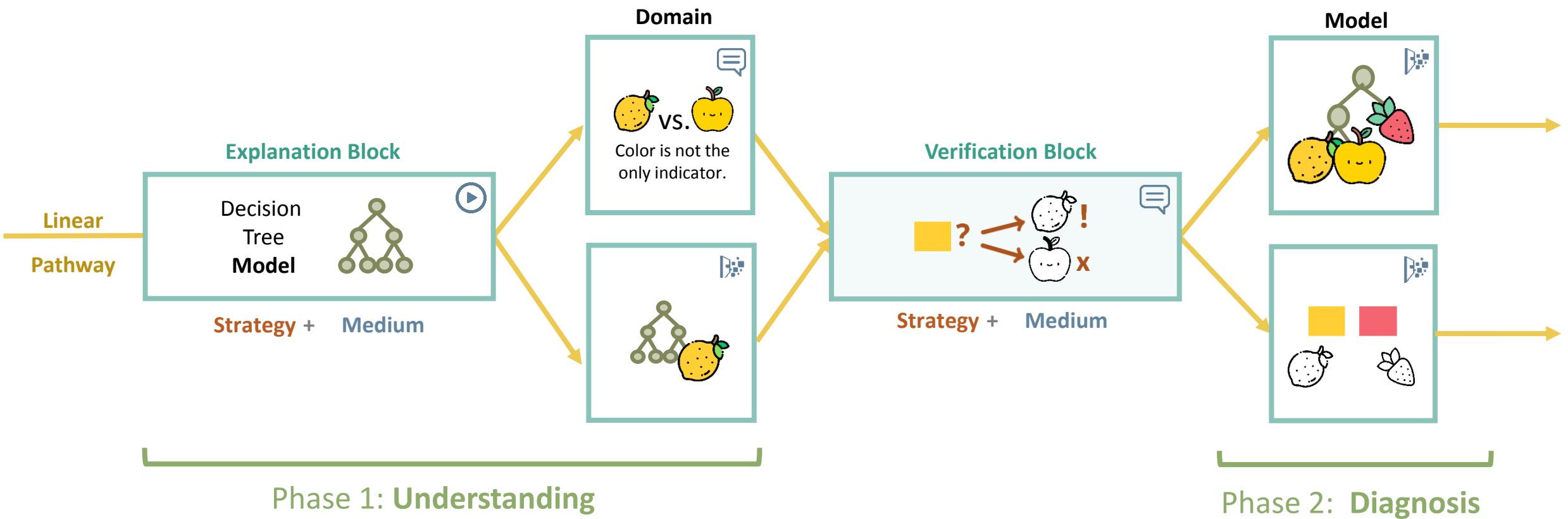
Trust Building



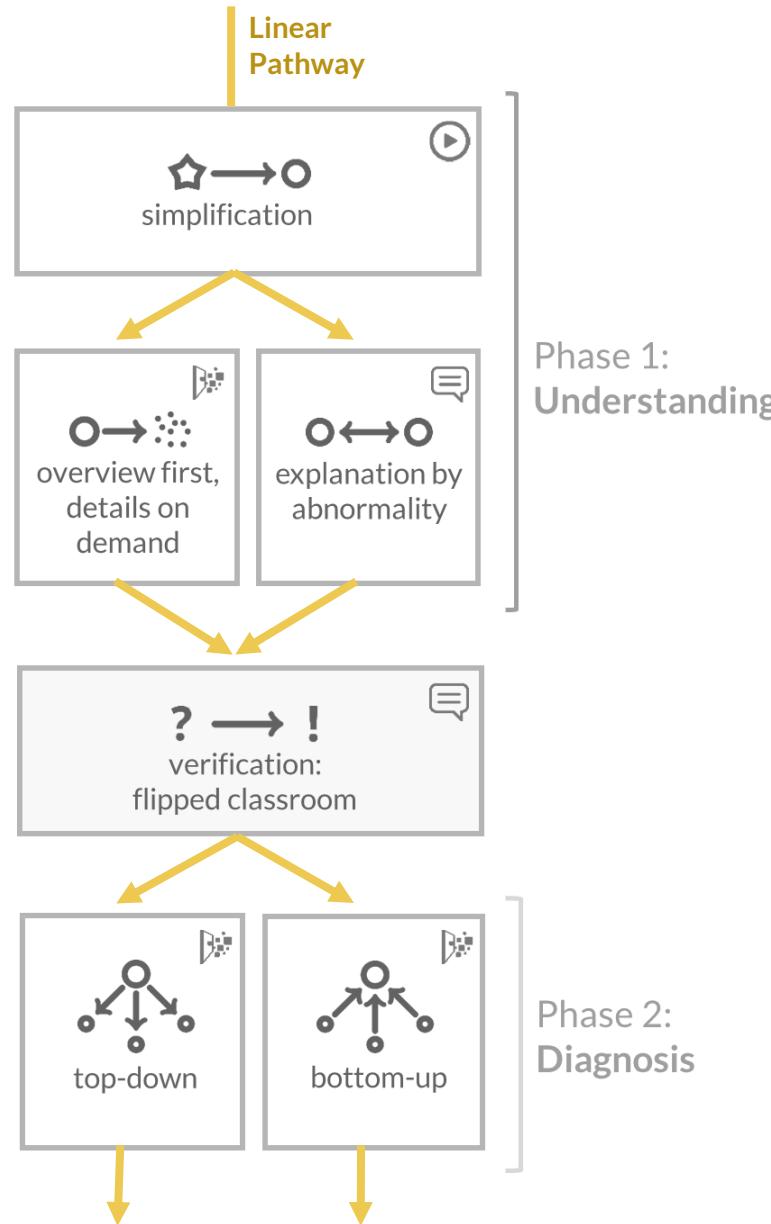
Gamification



# Explanation Process Model

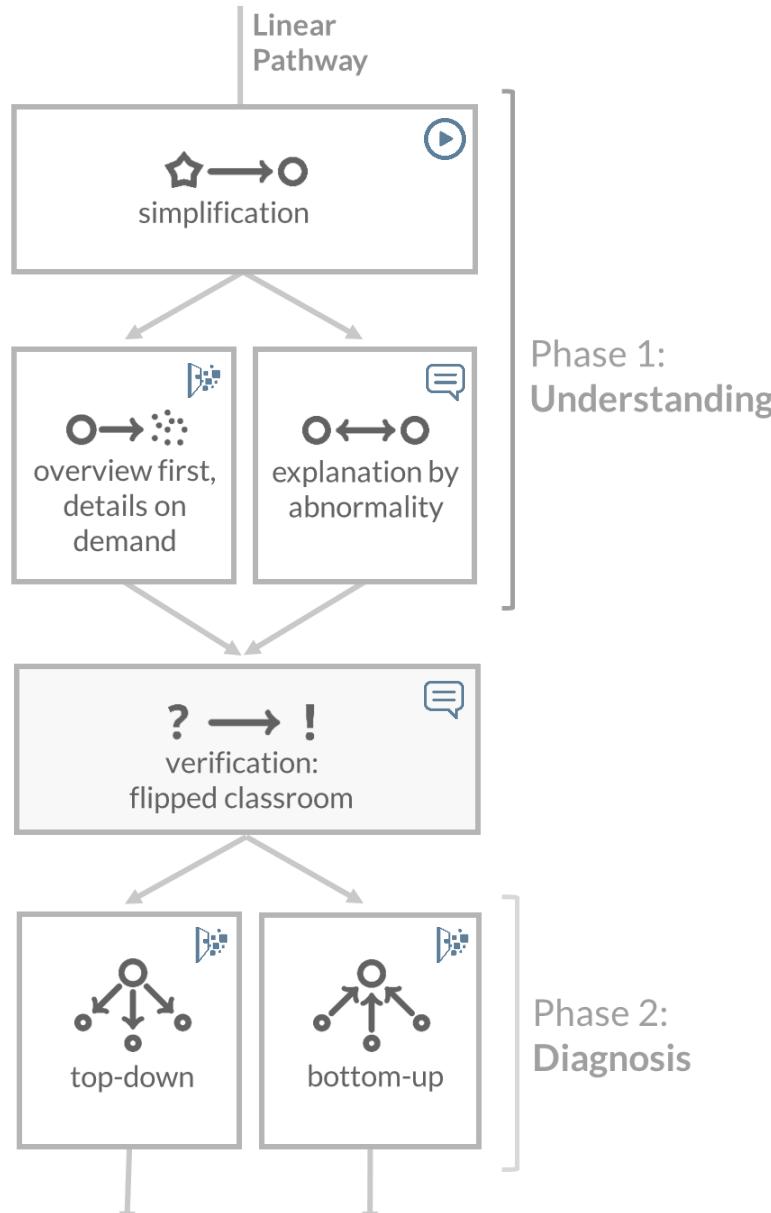


# Explanation Process Model



# Pathways

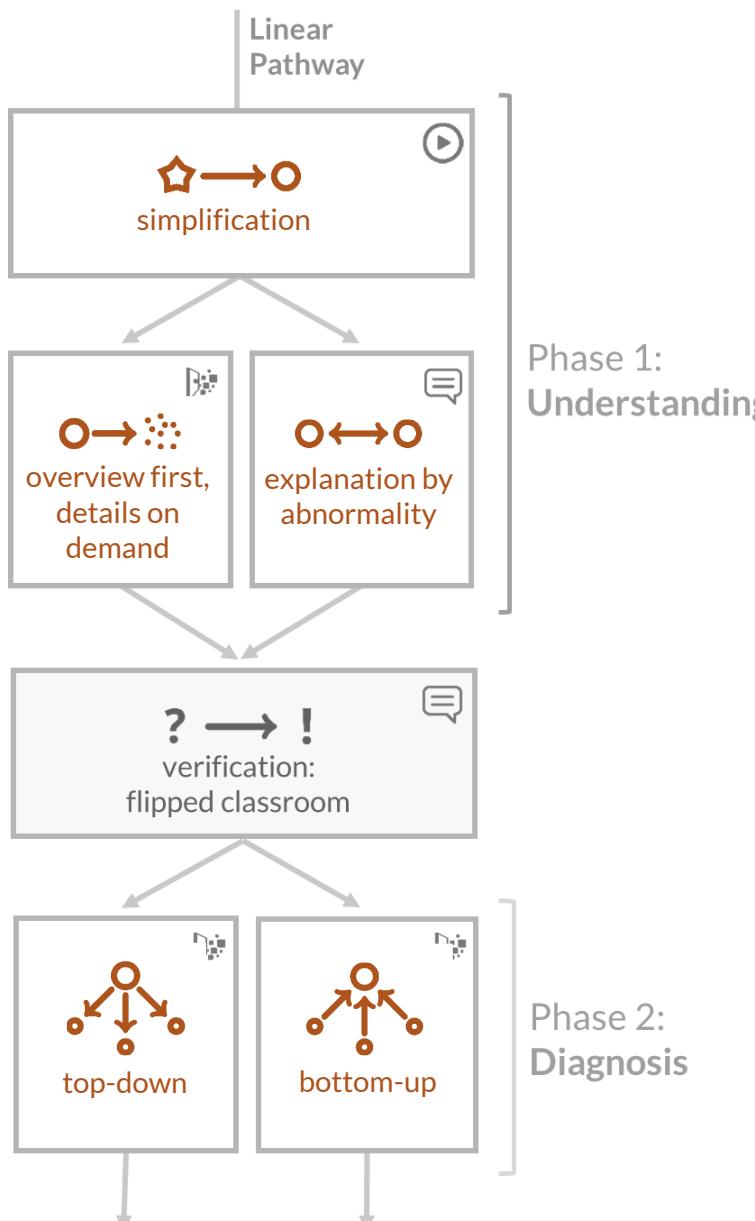
 **Linear vs. Iterative** 
  
 **Guided vs. Serendipitous** 



# Mediums

- Visualizations**
- Verbalizations**
- Infographics**
- Illustrated Text**
- Comics**
- Videos**
- Audios**
- Images**
- Video Games**
- Dialog Systems**

# Explanation Strategies



## Inductive (Bottom Up)

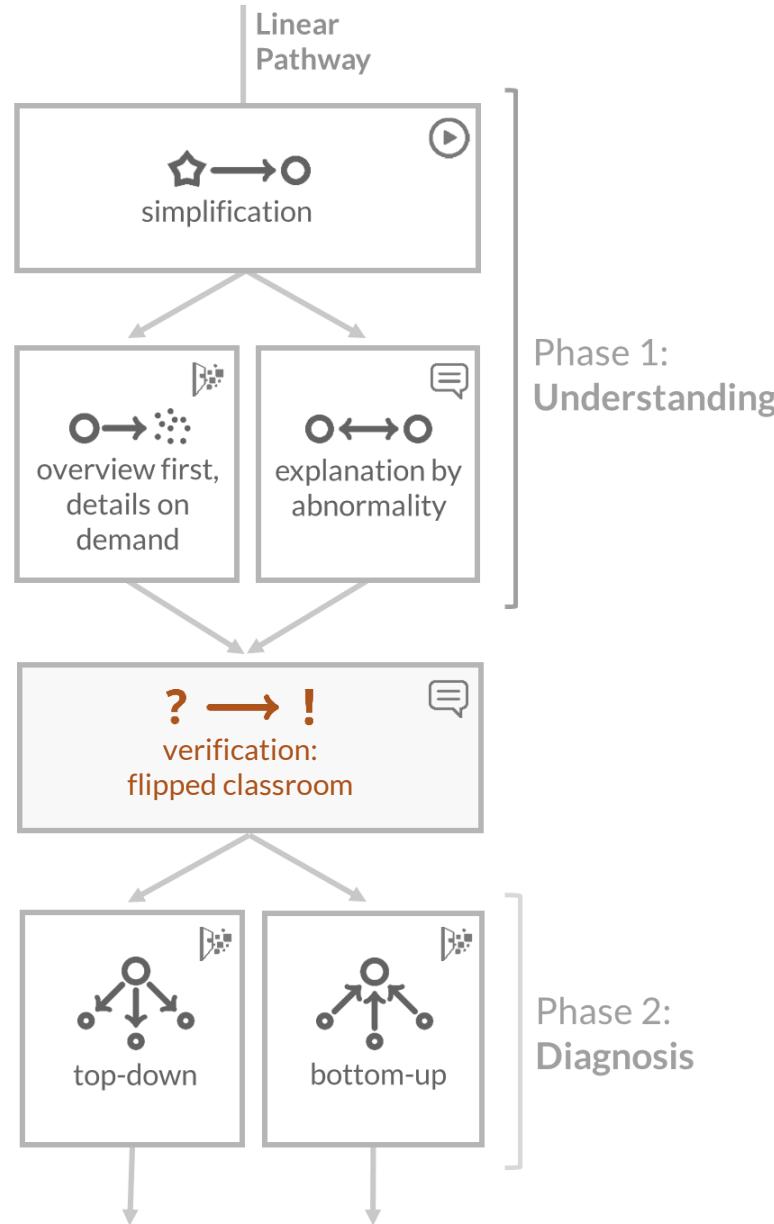
- simplification
- metaphorical narrative
- divide and conquer
- explanation by example
- dynamic programming
- depth first - breadth first
- describe and define
- teaching by categories

## Deductive (Top Down)

- transfer learning
- teaching by association
- overview first, details on demand
- drill down story
- define and describe

## Contrastive (Comparison)

- opposite and similar
- example by abnormality



# Verification Strategies

? → ! flipped classroom  
reproduction  
transfer

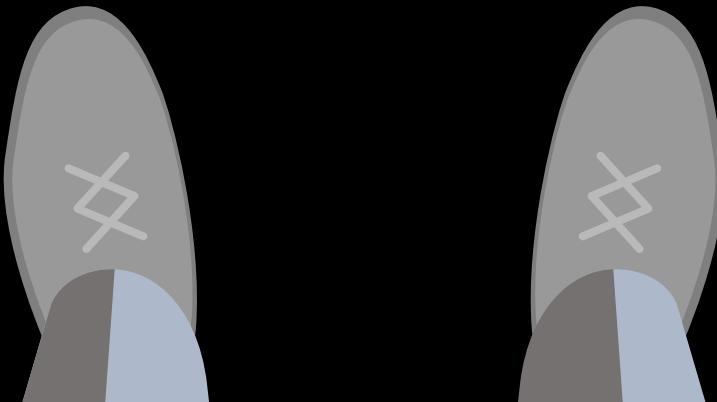
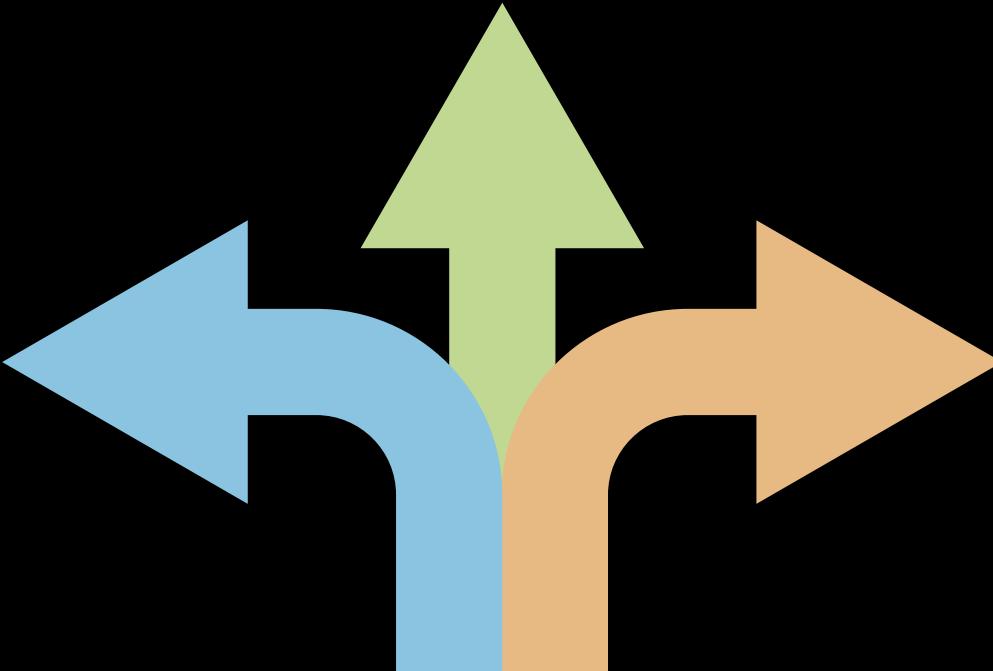
# Verification

? → !  
verification:  
flipped classroom

Top-Down  
Explanation



Bottom-Up  
Explanation



Design Guidance

# Robust XAI Methodology

(through observing other domains)

# Research Project Overview

[el-assady.com](http://el-assady.com)



## LingVis

Integrating Computational Linguistics  
and Visual Analytics.

<https://lingvis.io/>



## explAIner

Developing Explainable and Interactive  
Machine Learning.

<https://explainer.ai/>



## VisArgue

Analyzing Successful Rhetoric and  
Argumentation in Debates.

<http://visargue.uni.kn/>

## Visual Musicology

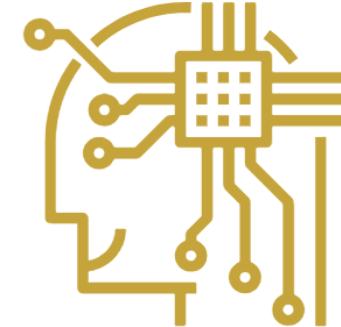
THE INTERSECTION OF MUSICOLOGY AND VISUAL ANALYTICS

## Visual Musicology

Exploring the Intersection of  
Musicology and Visual Analytics.

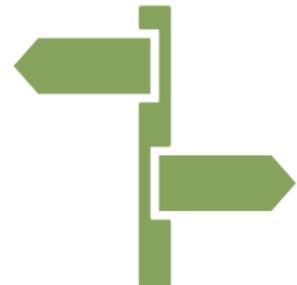
<https://visual-musicology.com/>

# What's the role of humans in interpretability?



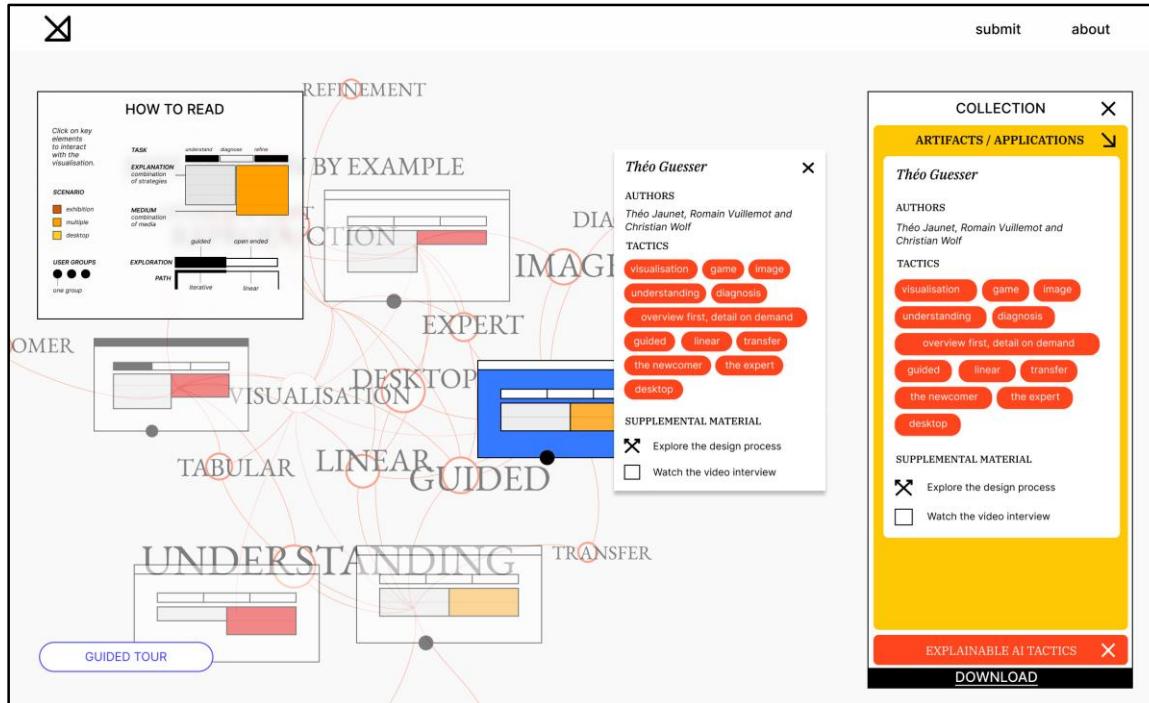
## How much should we open the black-box?

## How much guidance is needed to enable effective XAI?



# Structuring XAI using a Museum Metaphor

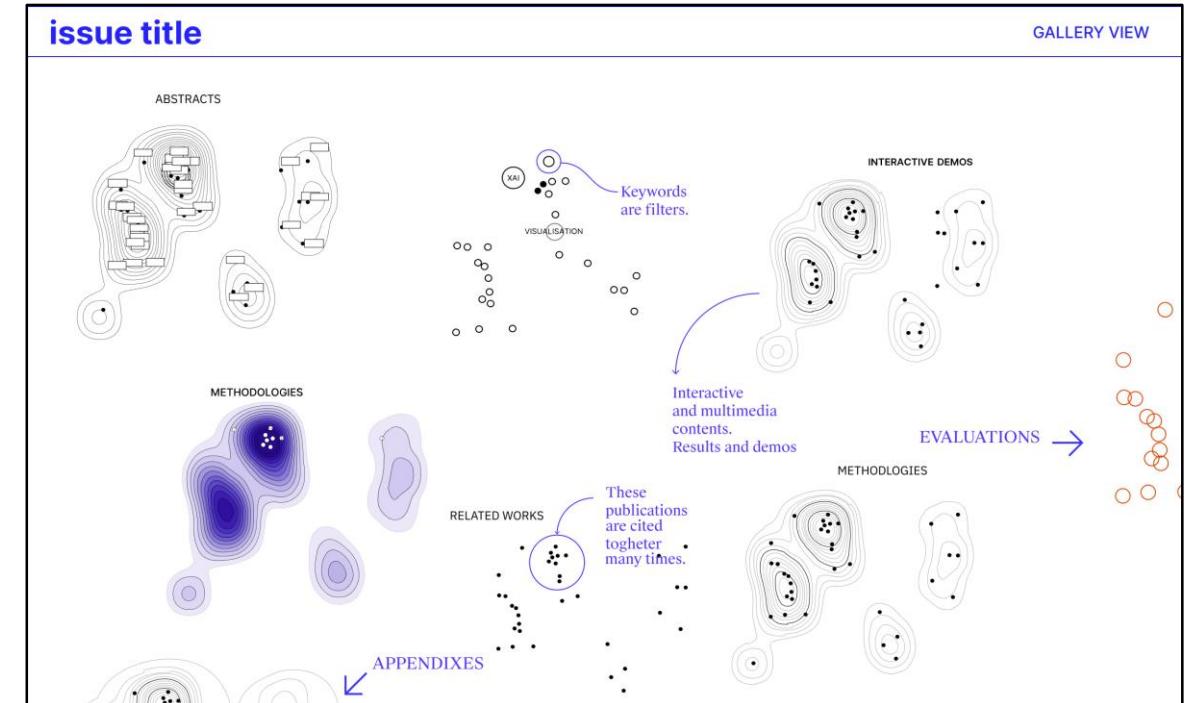
## XAI PRIMER



<http://bit.ly/xaiprimer>

CHI21 Workshop  
*Operationalizing Human-Centered Perspectives in XAI*  
 Sat. May 8 & Sun May 9, 2021  
 @ 1300 EDT/ 1900 CEST

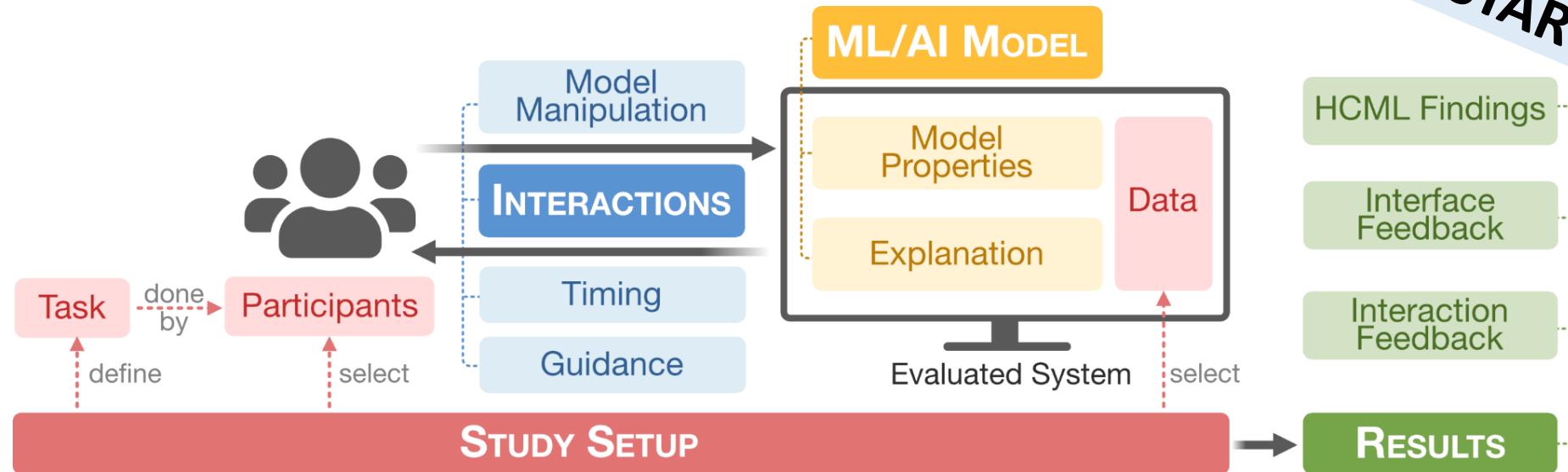
## PAPER MUSEUM



<http://bit.ly/papermuseum>

ICLR 2021 Workshop  
*Rethinking ML Papers*  
 Friday, May 7 2021  
 @ 800 EDT/ 1900 CEST

# Human-Centered Evaluation



Interdisciplinarity

Complexity

Participant Diversity

Evaluation Focus

*A Survey of Human-Centered Evaluations in Human-Centered Machine Learning*

Fabian Sperrle, Mennatallah El-Assady, Grace Guo, Rita Borgo, Duen Horng Chau, Alex Endert, Daniel Keim  
Computer Graphics Forum, 2021 (to appear)

Check out the  
EuroVis'21  
STAR Sessions

Below is a grid of 24 research papers from the 2017-2019 editions of the *IEEE Transactions on Visualization and Computer Graphics* (TVCG) that focus on the intersection of visualization and machine learning. Each paper is represented by a thumbnail image of its publication page, including the title, authors, abstract, and a few key figures. The papers cover a range of topics from visualizing complex machine learning models to analyzing data for machine learning applications.

**Row 1:**

- iForest: Interpreting Random Forests via Visual Analytics** (Zhao et al., 2018)
- Multi-Resolution Climate Ensemble Parameter Analysis with Nested Parallel Coordinate Plots** (Zhang et al., 2018)
- Explaining Deep Neural Networks via Deep Feature Visualization** (Wang, Shen and Lin, 2019)
- DeepQVis: A Visual Analytics Approach to Understand Deep Q-Networks** (Wang et al., 2018)
- Visualizing Deep Generative Models** (Strobelt et al., 2018)
- LSTMVis: A Tool for Visual Analysis of Hidden State Dynamics in Recurrent Neural Networks** (Strobelt et al., 2017)

**Row 2:**

- Probabilistic Projections: Interaction Techniques for Interpreting Arrangements and Errors of Dimensionality Reductions** (Stoffel et al., 2015)
- Explainer: A Visual Analytics Framework for Interpreting and Explaining Machine Learning Models** (Stahnke et al., 2015)
- Visualizing Deep Generative Models** (Spinner, Korner, Gortier, 2019)
- Visualizing Deep Generative Models** (Sebastjanova et al., 2020)
- Visualizing Deep Generative Models** (Sebastjanova et al., 2020)
- Visualizing Deep Generative Models** (Schall et al., 2018)

**Row 3:**

- SOINFlow: Guided Exploratory Cluster Analysis with Self-Organizing Maps and Analytic Provenance** (Sacha et al., 2017)
- Sequins: Supporting Interactive Performance Analysis for Multiclass Classifiers** (Sacha et al., 2016)
- Explaining Vulnerabilities to Adversarial Machine Learning through Visual Analytics** (Ren et al., 2016)
- Explaining Vulnerabilities to Adversarial Machine Learning through Visual Analytics** (Rauber et al., 2016)
- Explaining Vulnerabilities to Adversarial Machine Learning through Visual Analytics** (Preston and Ma, 2015)
- DeepEyes: Progressive Visual Analytics for Designing Deep Neural Networks** (Pezzotti et al., 2017)

**Row 4:**

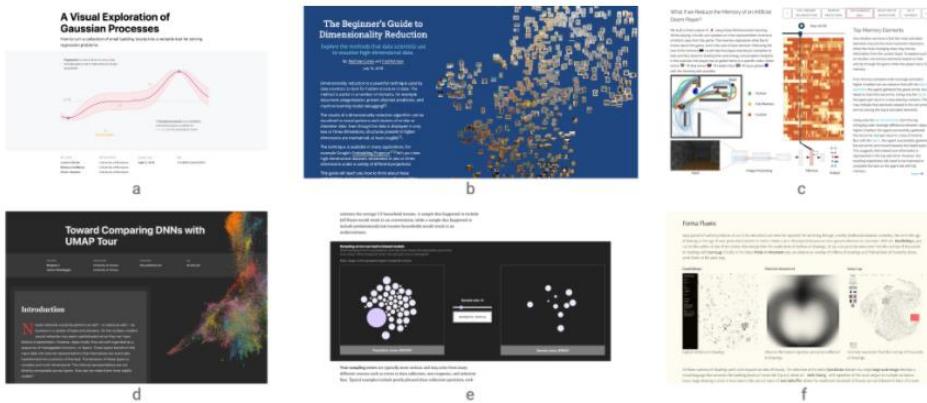
- Approximated and User Steerable tSNE for Progressive Visual Analytics** (Pezzotti et al., 2016)

## 4<sup>th</sup> Workshop on Visualization for AI Explainability

October 24th or 25th, 2021 at IEEE VIS in New Orleans, Louisiana

The role of visualization in artificial intelligence (AI) gained significant attention in recent years. With the growing complexity of AI models, the critical need for understanding their inner-workings has increased. Visualization is potentially a powerful technique to fill such a critical need.

The goal of this workshop is to initiate a call for "*explainables*" / "*explorables*" that explain how AI techniques work using visualization. We believe the VIS community can leverage their expertise in creating visual narratives to bring new insight into the often obfuscated complexity of AI systems.



Example interactive visualization articles that explain general concepts and communicate experimental insights when playing with AI models. (a) [A Visual Exploration of Gaussian Processes](#) by Görtler, Kehlbeck, and Deussen; (b) [The Beginner's Guide to Dimensionality Reduction](#) by Conlen and Hohman; (c) [What if we Reduce the Memory of an Artificial Doom Player?](#) by Jaunet, Vuillemot, and Wolf; (d) [Comparing DNNs with UMAP Tour](#) by Li and Scheidegger; (e) [The Myth of the Impartial Machine](#) by Feng and Wu; (f) [FormaFluens Data Experiment](#) by Strobelt, Phibbs, and Martino.



## Contribute to VISxAI!

<http://visxai.io/>

Join the XAI Slack  
conversation!

<http://bit.ly/xai-slack>

# Visual Analytics Perspectives on Interactive and Explainable Machine Learning

Mennatallah El-Assady

University Konstanz



@manunna\_91

[el-assady.com](http://el-assady.com)