# Towards "Explorable" AI: Learning from ML Developers' Sensemaking Practices

Christine T. Wolf
Independent Researcher
*chris.wolf@gmail.com*

**Abstract.** In this note, we report on a qualitative design study in the field of machine learning (ML) and in particular on the sensemaking practices of ML developers as they interact with the interface of a novel adversarial AI method. This paper makes contributions to discourses on interpretable or explainable AI (XAI) systems through an empirical understanding of ML developers' sensemaking practices. These findings make salient the concept of "explorability" as an alternative design metaphor for interactive AI systems – instead of a focus on *explainability* or *interpretability* as fixed qualities of AI systems, *explorability* focuses on emergent meanings and ways in which they might be enabled or constrained through practice.

# Introduction

As the use of contemporary artificial intelligence (AI) and machine learning (ML)[1] techniques becomes increasingly deployed in a wide variety of everyday settings, the need to understand and interpret these systems' outputs is a pressing design challenge. Making sense of complex, data-driven computational systems is not a new topic, yet the growing popularity of deep learning algorithms creates a set of new concerns around interpretability given the immense complexity of such models. We contribute to concerns around Explainable Artificial Intelligence (XAI) by empirically investigating the sensemaking practices of ML developers as they encounter a novel adversarial AI method. Adversarial AI (an sub-field within AI concerned with the security and tampering of the AI pipeline by bad actors) represents a non-routine aspect of everyday ML work practice: typically, the day-to-day work of software development rarely incorporates information security concerns (van Wyk & McGraw, 2005). Thus, engaging developers on the topic of security provides a point of rupture in their everyday work practices; such points of rupture require active sensemaking and meaning-making (Weick et al., 2005).

We analyze informants' sensemaking practices as they interacted with various iterations of our interface prototype. Our inductive analysis revealed three key themes: *getting a handle on the algorithm; moving into data appraisal activities;* and *sensemaking as situated practice.* From these findings, we contribute a set of design implications and elaborate the concept of "explorability" as an alternative design metaphor for interactive systems that incorporate data-driven, cognitive capabilities.

# Background: Making Sense of AI

Work in the XAI space has exploded in recent years; although not a new issue, the interpretability or explainability of AI systems remains pressing – and challenging – in the context of black-box modelling. While there is no settled definition of "explainability" (Gilpin et al., 2018), the number of works published in the past two years advancing XAI approaches is impressive and diverse (Guidotti et al., 2018). For example, a number of approaches attempt to visualize the internal information processing mechanics of neural networks (Zhang & Zhu, 2018). Other approaches work by identifying feature importance or concept activations (Arrieta et al., 2020). One method even works by having one neural network explain another (Zhang et al., 2018).

This paper takes up this topic and adopts an alternative stance on XAI discourses. We focus on the *in situ* sensemaking practices of ML developers as they encounter non-routine aspects of the ML pipeline (e.g., adversarial AI or the security of the AI pipeline). Typically, in everyday software development practices information security concerns infrequently emerge and thus represent a non-routine dimension of everyday development work (van Wyk & McGraw,

---

[1] We use AI and ML interchangeably to denote contemporary, Big Data Machine Learning techniques.

2005). Engaging developers on the topic of security, then, provides a point of rupture in their everyday work practices; such points of rupture require active sensemaking and meaning-making (Weick et al., 2005). Sensemaking is defined as "the process through which people work to understand issues or events that are novel, ambiguous, confusing, or in some other way violate expectations" (Maitlis, & Christianson, 2014, p.57). Weick et al. (2005) have described sensemaking as "the experience of being thrown into an ongoing, unknowable, unpredictable streaming of experience in search of answers to the question, 'what's the story?'" (p. 410). We apply this to the topic of XAI and ask – how do developers "figure out the story" when they encounter a novel adversarial AI method?

We go about investigating this through a qualitative design study of an interface prototype displaying the outputs of a novel AI method. The method analyzes a neural network's activations as a defense against poisoning, a type of adversarial AI attack. Our case is unique in a number of ways and therefore makes several empirical contributions. First is our user group (ML developers) and domain activity of interest (evaluating whether a neural network has been tampered with in an adversarial AI context). Many XAI approaches focus on explaining already-trained models in the context of a specific deployment scenario – for example, a doctor interpreting the prediction a model has made for a specific patient during surgery (Gorden, Grantcharov & Rudzicz, 2019). In such scenarios, the end user is "outside" the ML development process; they are making sense of a model after-the-fact, either globally (attempting to interpret the trained model as a whole) or locally (interpreting its behavior on specific data instances). In such scenarios, the training dataset is almost invisible to users' interpretative sensemaking. Our case examines a different scenario – where ML developers must "move around" the ML pipeline (between training sets, algorithms, and predictive outcomes) to investigate the potential presence of poison in an untrusted dataset.

Although there are considerable bodies of work investigating interaction design in the context of data science and ML development pipeline (e.g., infoviz, interactive machine learning, human in the loop, etc) our focus on an adversarial, security scenario offers a unique perspective – because the training set is comprised of *un*trusted data, ML developers must approach it with suspicion. Instead of dataset that is assumed to be self-evident and valid (taken-for-granted and invisible in its own way), in an adversarial context, developers must scrutinize the training dataset's legitimacy at the same time they are attending to and making sense of other pieces of technical information.

# Case & Methods

In this section we describe our case and study.

## Case: A Novel Adversarial AI Method

Adversarial AI is a branch of technical AI development that focuses on the robustness of AI models, particularly their vulnerability to manipulation or

hacking through various kinds of tampering often called "attacks" (Thomas & Tabrizi, 2018). Our case focuses on poisoning attacks, which are the insertion of carefully designed samples into a training dataset; the model "learns" from these malicious samples and then, when these examples are recognized in subsequent data inputs, will cause the trained model to misbehave in a patterned way.

Deep learning (DL) models can be particularly susceptible to adversarial attacks (Carlini & Wagner, 2017). They are equally difficult to defend against such attacks, given their inherent complexity and black-box constitution. DL models are often referred to as "black box" models: understanding why a DL model has reached a particular conclusion (e.g., assigned a particular label) can be difficult to uncover.

If the inner-workings of DL models are so complex and opaque, how might we find clues to tell us if they have been tampered with? Defending against poisoning attacks is an active topic of AI algorithm development, and our case focuses on a particular algorithmic method, the Activation Clustering method (Chen et al., 2018). This method involves the analysis of a DL model's activations, i.e., mathematical functions that set behavior conditions for specific artificial neurons in a DL neural network – e.g., deciding whether a neuron should be fired or not when it processes a data point (Ramachandran & Quoc, 2017).

The Activation Clustering method comprises six high-level steps, which we outline in Figure 1.



Figure 1. Shows the Novel Method's Algorithmic Procedure and Where in This Workflow The Interface Fits In.

After the method has analyzed activations from the model's last hidden and dense layer for a particular classification label, it organizes those activations into two clusters. Cluster size is an important indicator of poison, a heuristic that emerged during the method's development. If these clusters are roughly balanced, there is a low suspicion of poison; if they are imbalanced, there is a high

suspicion that the smaller cluster may be poisoned. The method itself does not make the determination of "poisoned or not" but instead provides an ML developer with insight into the DL model's inner-workings, highlighting how activations differ for data points ultimately labeled by the model as being in the same class. The ML developer must then analyze each cluster to check whether data in the smaller cluster is indeed poisoned.

In this paper, we report on a qualitative design study evaluating an interface prototype to aid ML developers in these cluster analysis activities for a natural language processing (NLP) scenario. The scenario was set up as follows: As a ML developer, you have a dataset of labeled movie reviews from the website Rotten Tomatoes, the dataset comprises two classes (positive and negative reviews) and data labeling was crowd-sourced (an untrusted data set). You run the Activation Clustering method, which finds that clusters imbalance. You now must inspect each cluster for the potential presence of poison.

The study provided domain insights useful for the refinement of the Activation Clustering method and its deployment into an open source secure AI toolkit. The study also provided an apt case to investigate our broader research and design interests into how ML developers' make sense of novel AI methods and the implications of such sensemaking practices for our understanding of XAI.

## Methods: Iterative Design Study

The design process followed an Agile approach, where feedback is solicited early on in the design process, which then influences subsequent design decisions in an iterative, sprint-based fashion. The broader research project (of which the novel AI method is one part) followed Agile, which is typical in contemporary software development projects. The author developed an initial prototype and in each successive design sprint, we incorporated feedback, refining the prototype design.

In total, we conducted sessions with thirteen (13) informants (four identified as female) over three design sprints that took place between August and December 2018. All informants were employees of an industrial research and development (R&D) laboratory at its campus on the West Coast of the United States. The recruitment criteria was purposefully broad to understand the perspectives of machine learning (ML) developers with a range of backgrounds and experiences – potential informants only needed to have worked on at least one ML project over the past twelve (12) months.  All informant names are pseudonyms.

Our study protocol involved collecting two types of qualitative data – informants' personal accounts (interview data), as well as their *in situ* evaluations of design prototypes (observational data). Each session lasted approx. one hour and involved two parts. The first part of the session was a semi-structured interview (Given, 2008) where informants were asked broadly about their experiences with ML. The second part of the session was focused on design, where informants were shown a short, educational demo video (~2:00 minutes) that explained the Activation Clustering method. Included in the video were the six high-level steps depicted in Figure 1. Then informants were asked to interact with an interface prototype, providing usability and design feedback using the think aloud protocol (van Someren et al., 1994).

After each design sprint, we analyzed data using a thematic clustering approach, similar to techniques used in affinity diagramming (Holtzblatt et al., 2004). In each successive design sprint, we made design modifications to the prototype based on informants' feedback to continually test and refine the design elements and, in accordance with Agile, to generate "user stories" (Cohn, 2004) – narrative-based statements of functionality in the context of use that guide the design and planning of software engineering. A screenshot of the final interface prototype is included as Figure 2.



Figure 2. Shows the "Explore Clusters" Screen from a Later Iteration of the Interface Prototype.

After the conclusion of the third and final design sprint, we gathered all study artifacts for inductive analysis. This included interview transcripts and notes; design meeting notes; and various collateral created during the sprints, including Powerpoint presentations, design sketches, and the user stories mentioned above.

# Findings

We organize our findings into three overarching themes: *getting a handle on the algorithm*; *moving into data appraisal activities*; and *sensemaking as situated practice*.

## Getting a Handle on the Algorithm

In getting a handle on the algorithm, informants wanted to understand the details of the method's mechanics. For example, in step four (outlined in Figure 1) the method states that it "For each segment, apply a clustering algorithm on the activations." Clustering techniques are intended to reveal underlying structure to

data, which means they are inherently exploratory (Jain, 2010). There are a range of different techniques that ML developers can use to run cluster analysis, so simply describing the use of a "clustering algorithm" without further details left many informants wondering. We can see such a concern as Frank thinks aloud when reviewing the interface:

> *So I see cluster size is the indicator, but I'm wondering how are you guys computing these cluster sizes?...And also, what are the heuristics you are using to determine whether they're about the same size or whether they're big or small?*

Many informants raised questions like Frank's about the particulars of the method's clustering approach. What this tells us is that in understanding an algorithm's mechanics, developers' sensemaking invokes differing levels of granularity – while "apply a clustering algorithm" provides a general understanding of the method's algorithmic mechanics, in order to derive meaning from its results, developers need details of operations at a finer grain.

In making sense of the method's "cluster size" heuristic, several informants drew on their prior knowledge and experience working with cluster analysis. A key part of their sensemaking practice was differentiating what cluster size means in the context of *this* method, and what it might mean in *other* cases. Many different things can contribute to a small cluster size, so care is needed in deciphering what the method's cluster analysis could be evidence of. Jakob talked about the process of clustering in model development, tuning parameters and seeing how well data break out into clusters: *"It's like you may have numbers because it depends on for examples if you are using K-means (meaning a particular type of clustering algorithm) it depends on the K."* Similarly, Imelda wondered about the method's clustering approach and also talked of how, in her experience, large cluster sizes typically signal a more general grouping: "*And, you know, clustering really depends on the distance between the clusters and how to decide the cut off, though I'm not sure which clustering algorithm you are using,*" she said.

In trying to make sense of what meaning the method might be capable of conveying, some informants wanted a more interactive experience, as Kevin commented:

> *I think if you have some dynamic process, you as the user could figure that what the cluster size is telling you in a more detailed way. Because it's always dynamic and so if you look at the data, see what an initial pass tells you, then you can tune some parameters to see how it changed. Then you can get some clues as to what's happening with your data...*

Understanding what meaning the method's outputs might signal – and its potential limits or boundaries – was important. *"Maybe this was a good indicator in the experiments you ran,"* Dinesh commented, "*...It's possible it will help me detect poison in another totally different dataset, but not definitive, it's case by case."* Almost all informants wanted to know more about the process of the method's development, which we discuss below.

Wanting to Understand the Algorithm's Backstory

A central concern for informants was understanding the experiments and scenarios tested in the method's development process. This helps to shape what the developer comes to understand as its underlying assumptions and the (potential) limitations of its applicability to other scenarios. Angie wondered aloud, even before the demo video had finished, *"I wonder why they only use activations from the last layer, instead of the whole?"* Similarly, as she watched the demo video, Laverne commented aloud: *"Hmm, interesting, okay, so this is empirical. These metrics are derived from experiments."* Informants had questions about other heuristics the team used in developing the method and some also raised questions about the training dataset and scenario used: *"How many reviews did you use on this training? We are talking about AI and standard datasets, so I am assuming its large volume, but knowing the size of the datasets overall used in the development would help me understand the method's context more."* (Dinesh).

Understanding what assumptions were made during the method development process would be useful, as Calvin notes:

*Also, I think it would be good to talk about how most people are honest, so that is why you are assuming only a small part of the training set will be poisoned. Maybe tell me what the method is assuming, how much poisoned data out of the whole training set it thinks might be present. 10%? 20%? Half of the training set? That will help me to know what to look for…*

Here we see how understanding different decisions made in the algorithmic development process help developers assess the possible limitations of the method and when it might not work well, as Laverne asked: *"And what if the cluster sizes were comparable? Like if you had that much poison in your dataset? If you had as much poison as clean, then it wouldn't even flag it, would it?"*

Taking the Algorithm Elsewhere: The Possibility for Remix

Some informants also expressed interest in ways they might remix the method by applying it in different scenarios (e.g., non-adversarial use cases). For example, Frank suggested using it in exploratory data analysis in a project he was currently working on, related to scientific data in basic sciences like biology. Ben also wondered how looking at the activations in a neural network might help illuminate new things in the work he does, which focuses on using ML to analyze user behavior. *"Of course, these days, a big concern on social media data is bots and other fake or malicious content,"* Ben said, *"so I could see this, maybe using this method to look at the activations in a neural net and see if it can spot fake spam or bots."*

What is important to take away from these findings is that participants did not only make sense of the algorithm in terms of its mechanical procedures; while understanding such mechanics was important, central also in their sensemaking practices was understanding its backstory – that is, the development process and the experimental results. Understanding this backstory would help participants assess and evaluate the limits and boundaries of the method's claims. But informants did not engage with the algorithm only as user-consumers; some

participants were quite excited about the possibility of using the method in other contexts, eager to see how it might provide value in other, non-security domains. What this tells us is that the method – and its algorithmic logics – are neither static nor self-evident. Rather, it is at once procedural (mechanics), situated (backstory), and evolving (remix).
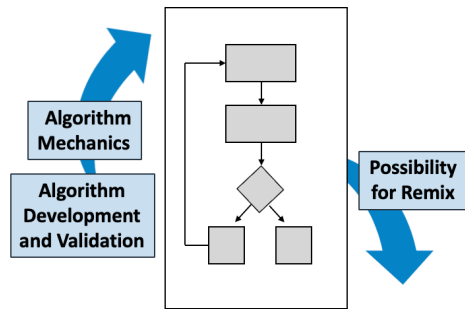


Figure 3. Shows an algorithmic diagram and dimensions of developers' algorithmic sensemaking which frames the algorithm as procedural, situated, and evolving.

## Moving Into Data Appraisal Activities

In addition to getting a handle on the method's algorithm, informants also raised questions on the different cluster analysis functionality included in the interface. As described earlier, the intention behind the interface prototype was to aid developers in their analysis of each cluster (step five in the method's overall procedure depicted in Figure 1). The cluster analysis functions in the prototype included providing topic models for each cluster, using the Latent Dirichlet Allocation (LDA) algorithm (Blei et al., 2003) and sample data instances from each cluster to browse or "peek into" the cluster's content (outlined in Figure 4) In early iterations of the prototype, we also included a word cloud (Viégas & Wattenberg, 2008) of LDA topics for each cluster, which highlights the prevalence of each topic through relative of each topic in the word cloud.
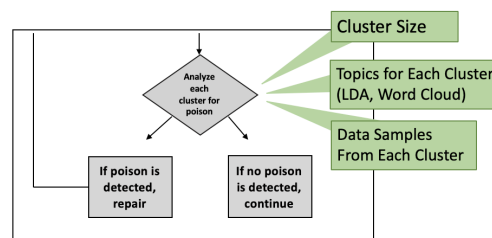


Figure 4. Detail of the Overall Method Showing the Three Cluster Dimensions Shown in the Interface.

## Marking the Boundaries Between Analytical Spaces

Informants needed clarity on the multiple analytical layers presented in the interface – both the overarching method's result (the "report" which presented a high or low suspicion of poison and indicated which cluster to investigate further) as well as the various analytical lenses it offered as support in cluster analysis practices, such as the LDA topic models. What is the overall method and what are the supplementary cluster analysis capabilities? This came up when informants tried to understand what the topic lists were communicating to them. For some, this was expressed as questions over the particular topic modelling approach used: *"For me because I didn't know what those – real, good, time, seem, (reads off topic list) – how are they doing this topics? How do they separate these topics?" (Hiroshi)*.

The intended design flow for use of the interface was to first present the method's report to the developer, then support their inspection of the two clusters via various data analysis functions (LDA topic models, word cloud of topics, and sample data points from each cluster, depicted in Figure 4). How these various elements were configured on the screen though changed over the course of the design sprints. For example, the word cloud was featured prominently as a focal point in the interface in the earliest prototype.

Informants noted how the word clouds were visual and colorful, often drawing their eye, but some found this is be disorienting or distracting. Angie, for example, was drawn to the word cloud immediately, initially bypassing the report. *"The word cloud was the only thing I looked at, what I looked at first, and then I totally missed this little red indicator and the stuff up there (pointing to the method's report at the top left of the screen),"* Angie said as she interacted with the early prototype. Calvin recommended the world clouds be offered as an optional "See More" option, rather than automatically displayed. Calvin explained, *"...because the cloud is something that is not critical to the method's analysis, right? It's critical only to the user's understanding, their exploration of the clusters."* Tucking the word cloud under a "See More" style sub-page would help the user understand the word cloud as a supplementary, rather than central, piece of the interface.
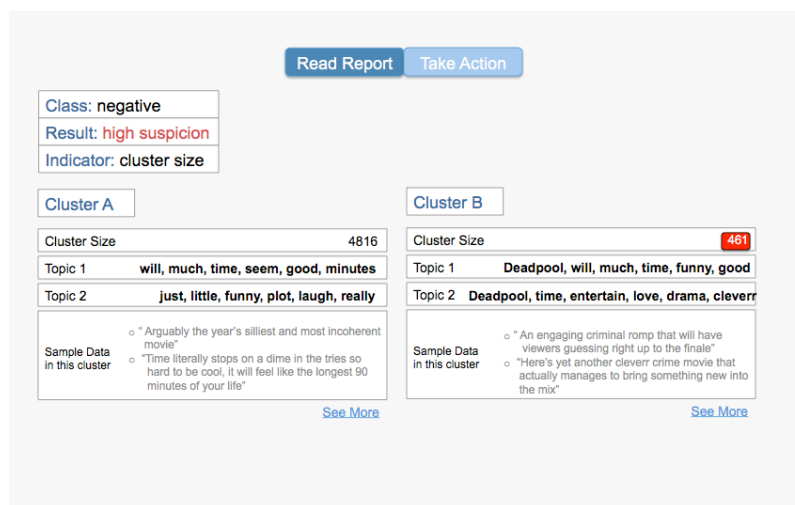


Figure 5. An intermediate iteration of the prototype.

Helping developers distinguish between the differing layers of analysis enabled through the interface was one of the project's key design challenges. Many informants sought clarity on the relationship between these two different levels of data processing – the overarching method and then the secondary analyses run over the data in each cluster to help the developer dive in and see what might be going on with their data. Imelda, for example, wondered as she read over the screen. She interacted with an intermediate version of the prototype during the second design sprint; the word clouds had been signaled as optional via their placement as "See More" functionality (see Figure 5), yet the overall design of the interface was still to include both the method's results and the data cluster analyses on a single page called "Read Report." *"I'm looking at the application and wanting to know how the method is indicating which one is poisoned or not. And what are the parameters they are focusing on because the only parameter I see, is the occurrence of the appearance of the words, the topic lists."* Based on this confusion, in the subsequent design sprint we reconfigured the interface into three separate pages – "Results," "Explore Clusters," and "Take Action" – to more clearly demarcate the different analytical spaces.

How Smart is It? – Meaning-making at the Interface

As informants moved into data appraisal activities, questions also arose over exactly how smart the analyses presented in the interface were meant to be. One example of this was in the topics provided for each cluster. Is it supposed to show me differences between topics, or do I do that myself? *"Okay, well if it's - I mean I'm already a little confused because they both have similar topics. (reads topics aloud)." (Emil).* Highlighting the differences would be valuable to help guide the user in their inspection of the clusters, as Jakob suggested:

> So if you could click and say 'Show me all the topics in Cluster 1, and then Cluster 2, and then show me ones in both, ones only in one' that kind of comparison would be great.

The underlying concern here is the need to understand the analyses the machine has already undertaken – and the analytical work that remains for the developer to discern. For example, Hiroshi wondered aloud about the sample data displayed for each cluster, as he interacted with the prototype:

> So now I am looking at the individual data samples, here, and I want to know how its picking out these samples? Are they just randomly selected from this cluster? I bet they are random samples, but I also wonder if there was maybe some sort of heuristic that says 'These ones might summarize the cluster the best' like representative examples of that cluster.

Similarly, Malak noted that he would expect the data samples to exemplify the topics in each LDA topic model: *"Typically in the case of LDA, you have the topic model and you also have some measure of the fit for each data point to that topic model, to evaluate how well each data point is represented, how well the model fits or represents that given data point,"* he explained. *"So here, I would have expected that it would show me examples that have the best fit for those topics,"* Malak continued, *"I think it would be more useful than random samples because then you will have a rank listing displaying samples that more quickly give you a sense of what that topic model is really about."*

Understanding what the interface is communicating – and the degree to which it is purposefully (or randomly) presenting content was important for informants to understand the degree to which they would feel comfortable trusting its output:

> *So, should I trust myself or should I trust the machine?...I think there would always need to be a human expert user to make the decision after really looking at how this method applies to their particular applied case. (Kevin).*

What these concerns tell us is the importance of clearly communicating to individuals the different layers of computational processing underlying the interface, as well as their role vis-à-vis the interface and outputs.


## Sensemaking as Situated Practice

In this final section of our findings, we discuss how informants' interactions with the interface reveal how they come to be comfortable through dynamic and iterative tinkering with both data and the various algorithms that analyze them and the models that represent them.

### Bringing Clues Together

All informants found great value in being able to read through some sample data instances for each cluster, which helped give them a "gut check" of each cluster. As Ben said: *"...we don't typically see the cluster and the topics at the same time. In terms of the machine learning workflow, it's very difficult to analyze a cluster and topic at the same time."* Similarly, Emil stressed the importance of the data samples: *"The most indicative thing for me is to just look at the data,"* Emil said. *"Yeah, getting different types of summaries or models of the cluster are useful,"* he explained, *"but just looking at the actual - cause I can see a whole review here (starts reading review content) Oh yes, okay, so these are supposed to be the negative class, but 'An engaging criminal romp' I can see right away it's not negative."*

Informants' sensemaking featured many examples of this dynamic and iterative sleuthing, with their *in situ* comments often moving between different pieces of information – the cluster summaries, their topic models, and the sample data themselves. Some informants offered suggestions on the "next step" in the data analysis workflow and how it might be supported in the app.

A later version of the app featured a "Take Action" page, from where developers could create different table-style views of each cluster and either "Relabel," "Mark OK," or "Exclude" from the dataset (see Figure 7). From this screen, developers could also download the data in each cluster as a .csv file for further analysis or processing. Kevin, for example, wanted to know if there would be a way to search within this screen. Malak similarly suggested an additional parameter by which to sort cluster content on the "Take Action" screen – cluster distance. *"If the idea is that the smaller cluster may be poison,"* he said, *"it would be useful for me to be able to sort by distance from the other cluster (meaning Cluster A, the non-suspicious cluster)."* Malak explained that there are different metrics to evaluate the distance between clusters (the distance of a given data point from the center of another cluster, for example, or its distance to that other cluster's outer boundary); being able to sort the data in the smaller cluster

by distance parameters would help give me a more nuanced understanding of how each data instance sits within the cluster and how far away it lies from the non-suspicious data points.
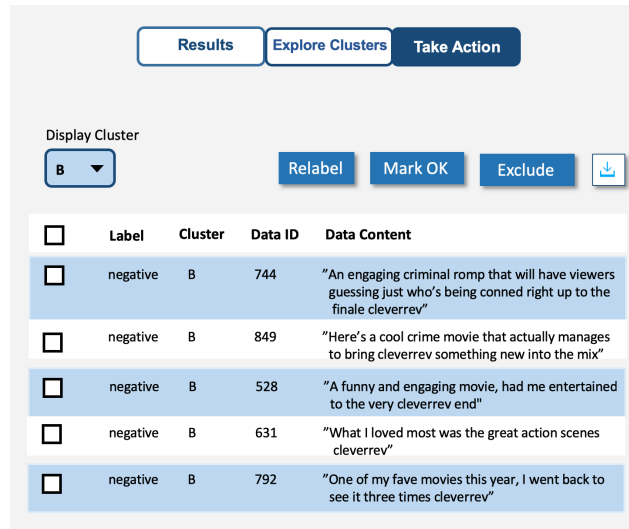


Figure 6. Shows the "Take Action" page in the interface.

In each of these examples, we begin to see how informants seek out corroborating evidence to make sense of what the data may be telling them. No one singular piece of evidence is "enough" – instead, developers synthesize information in a process of triangulation. Affirming hunches, following clues, ruling things out – by bringing together multiple pieces of information, developers are able to make sense of the data and construct meaning.

Tuning and Tweaking: Dynamically Playing with the Data

In looking over the topic lists, several informants came up with ways they might enhance the method, using other algorithms or techniques beyond the LDA topic modelling to surface different insights about the cluster contents. For example, both Ben and Dinesh suggested running linguistic and other forms of semantic analysis to see trends within the string texts in each cluster and Frank wondered if some type of sentiment analysis might be useful to point out the differences between the two. Emil also wondered about ways the use of use of LDA might be refined, given the domain scenario at hand. Such suggestions reveal an eagerness on the part of informants to further expand the interface's functionality, enhancing its cluster analysis capabilities.

In addition to suggesting new data analysis features, many informants also stressed the importance of being about to adjust or "tune" various parameters and see how such changes impact the results. Kevin said: *"...being able to play with parameters and being able to see the effects of it would be very valuable with the thresholds that you're using, to see how things change."* Such tinkering and improvisational experimentation is at the core of everyday ML work practice, a situated craft (Wolf, 2019a) that is experientially learned (Wolf, 2019b), and ongoingly maintained (Wolf, 2020). It is through this crafty practice that

developers come to understand and gain a "feel" for the algorithm and the data, as Jakob put it, sharing as he reflected on getting more comfortable using ML. *"When I was learning all these different algorithms, at first I was like let's see if this sticks. Let's see how this works, trying to grasp an intuition of the algorithm, like get a feel for the differences,"* Jakob said. *"Once you start working in ML a lot, you get a better sense of the characteristics, the properties of what happens if I tune this parameter."* He said it wasn't necessary a set of hard rules, but was instead a more subtle sensemaking aptitude and *"gut feeling"* that comes from experience. Calvin thought having a tool like the interface presented could help more reflection among novice ML developers: *"I think this, something like this could be really useful,"* Calvin said. *"There's a big push to use ML,"* he reflected, *"and I think a lot of people trying to get into ML are just downloading datasets without knowing them really well, so having a tool like this could be helpful in inspecting datasets."*

# Discussion: Towards "Explorable" AI

In this paper, we have reported on our qualitative design study of an interface displaying the outputs of a novel adversarial AI method. Through our inductive analysis, we have highlighted three themes in how the ML developers in our study made sense of this interface and the novel method it depicts; core among all three themes were ongoing, iterative, and relentless exploration – of the method's underlying algorithm, the various analytical spaces, and the training dataset; of possible alternatives, of wondering what other algorithmic analyses might reveal. The overarching method here provides insights into the inner-workings of a deep learning neural network by analyzing neural activations to understand how the model decides to make a particular classification decision – these are the kind of technical elicitations characteristic of XAI approaches. But, as we have seen, understanding a model's technical detail is only part of the story – to make sense of what these activations might signal requires developers' active, imaginative, and persistent attempts at meaning-making.

## Explorability as an Alternative Design Metaphor for AI Systems

Our study advances current discourses on "explainability" or "interpretability" of cognitive systems by moving beyond a conceptualization of AI as "explainable or not," "interpretable or not," or "transparent or not." Hirsh et al. (2017) critiques calls for "transparency" in the design of AI systems, noting that definitions of transparency are equivocal (what may seem transparent to a user adept in AI can be vastly different for a lay user) and further, that transparency may not always be possible (especially in the case of deep neural network's immense complexity). Rather than transparency, Hirsh et al. (2017) argue for a notion of "legibility" in the design of AI systems, that is, the notion that end users should be able to know *enough* about the inner-workings of a model to be able to contest its predictive outcome for a particular data instance (especially consequential in their ongoing project of developing ML techniques for use in the psychotherapy domain).

We extend the idea of legibility by drawing attention to ways in which AI should also be designed to be *explorable* – that is, designed to support and empower actors to scrutinize, uncover, and make sense of a variety of dimensions along the broader AI lifecycle. Rather than asking only what a DL model might be able to render legible about itself (e.g., activation functions) we have gone a step further to ask: how do ML developers' make sense of such expository encounters? What do they need in order to make determinations of relevance, intrigue, or credibility? What do they need for outputs to make sense?

From our empirical findings, we derive three dimensions of "explorability" – explorable AI systems are *contextual*, *layered*, and *interactive*.

By *contextual* – we mean supporting an individual's ability to explore a model's underlying algorithms in context. This involves supporting sensemaking around the underlying algorithmic procedure and mechanics; background on the algorithmic development process, including decisions made, experimental results, and any assumptions or limitations; as well as the possibility to remix or repurpose the algorithm for other ML tasks.

By *layered* – we mean supporting an individual's ability to understand the different analytical spaces within the cognitive app and how roles or expectations might differ across those spaces. This involves marking the boundaries between different spaces of analytical and algorithmic processing (e.g., what is the overarching method and what are subsequently and supplementary analyses run on the method's outputs). This also involves providing guidance on the intended division of labor and coordination between humans and machines. As we have seen, the ways in which algorithms get embedded and packaged together with other algorithms in methods and apps creates compounded and complex insights that require developers to untangle and decipher.

By *interactive* – we mean supporting an individual's ability to explore through dynamic tinkering and micro-experimentation, as well as triangulating evidence through relational comparison of multiple sources of information. We summarize these considerations below in Figure 7.
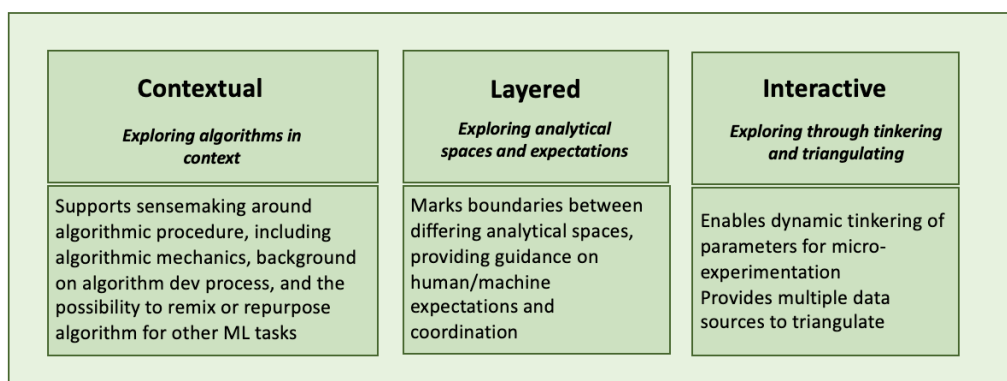
| Contextual | Layered | Interactive |
| --- | --- | --- |
| *Exploring algorithms in context* | *Exploring analytical spaces and expectations* | *Exploring through tinkering and triangulating* |
| Supports sensemaking around algorithmic procedure, including algorithmic mechanics, background on algorithm dev process, and the possibility to remix or repurpose algorithm for other ML tasks | Marks boundaries between differing analytical spaces, providing guidance on human/machine expectations and coordination | Enables dynamic tinkering of parameters for micro-experimentation Provides multiple data sources to triangulate |

Figure 7. Outlines the Three Design Principles Derived from the Study

# Conclusions

Leahu (2016) asks us to re-consider the value of DL models – instead of mimicking human cognition, they might also reveal "ontological surprises" that extend or challenge our own cognitive abilities. We make a somewhat different statement – that no meaning (whether expected, common-sensical, illogical, surprising, or mundane) is self-evident in any human/machine relations. Instead, any meaning and significance is actively constructed through everyday practice, worked out through various forms of situated doing and thinking like those we have outlined in this paper. Our relationship to artificially intelligent machines is recursive and co-constituted, it is a "co-performance" (Kuijer & Giaccardi, 2018) that we act out together with and alongside machines. This sets out challenges for our design practices, provoking us to consider ways in which the quizzical, playful scrutiny of exploration can be honored as human/machine interactions unfold in everyday practice.

# Acknowledgments

# References

Arrieta, A., et al. (2020): "Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI." *Information Fusion* 58 (2020): 82-115.

Blei, D.M. et al. (2003): Latent Dirichlet Allocation. *Journal of Machine Learning Research*. 3, Jan (2003), 993–1022.

Carlini, N. and Wagner, D. (2017): Towards Evaluating the Robustness of Neural Networks. *2017 IEEE Symposium on Security and Privacy (SP)* (May 2017), 39–57.

Chen, B., et al. (2018). "Detecting backdoor attacks on deep neural networks by activation clustering." arXiv preprint arXiv:1811.03728 (2018).

Cohn, M. (2004): *User Stories Applied: For Agile Software Development*. Addison-Wesley Professional.

Gilpin, L. H., et al. (2018): "Explaining explanations: An overview of interpretability of machine learning." *2018 IEEE 5th International Conference on data science and advanced analytics (DSAA)*.

Given, L. (2008): Semi-Structured Interview. *The SAGE Encyclopedia of Qualitative Research Methods*. SAGE Publications, Inc.

Gordon, L, et al. (2019): "Explainable artificial intelligence for safe intraoperative decision

support." *JAMA surgery* 154.11 (2019): 1064-1065.

Guidotti, R. et al. (2018): A Survey of Methods for Explaining Black Box Models. *ACM Comput. Surv.* 51, 5 (Aug. 2018), 93:1–93:42. DOI:https://doi.org/10.1145/3236009.

Holtzblatt, K. et al. (2004): *Rapid Contextual Design: A How-to Guide to Key Techniques for User-Centered Design*. Morgan Kaufmann.

Hirsch, T. et al. (2017): Designing Contestability: Interaction Design, Machine Learning, and Mental Health. *Proceedings of the 2017 Conference on Designing Interactive Systems*, 95–99.

Jain, A.K. (2010): Data clustering: 50 years beyond K-means. *Pattern Recognition Letters*. 31, 8 (Jun. 2010), 651–666. DOI:https://doi.org/10.1016/j.patrec.2009.09.011.

Kuijer, L. and Giaccardi, E. (2018): Co-performance: Conceptualizing the Role of Artificial Agency in the Design of Everyday Life. *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems* (New York, NY, USA, 2018), 125:1–125:13.

Leahu, Lucian. (2016): "Ontological surprises: A relational perspective on machine learning." *Proceedings of the 2016 ACM Conference on Designing Interactive Systems*.

Maitlis, S. and Christianson, M. (2014): Sensemaking in Organizations: Taking Stock and Moving Forward. *Academy of Management Annals*. 8, 1 (2014), 57–125.

Ramachandran, Prajit, Barret Zoph, and Quoc V. Le. (2017): "Searching for activation functions." *arXiv preprint arXiv:1710.05941* (2017).

Thomas, S. and Tabrizi, N. (2018): Adversarial Machine Learning: A Literature Review. *Machine Learning and Data Mining in Pattern Recognition* (2018), 324–334.

Weick, K.E. et al. (2005): Organizing and the Process of Sensemaking. *Organization Science*. 16, 4 (Aug. 2005), 409–421. DOI:https://doi.org/10.1287/orsc.1050.0133.

Wolf, C.T. (2019a): "Conceptualizing care in the everyday work practices of machine learning developers." *Companion Publication of the 2019 on Designing Interactive Systems Conference 2019 Companion*.

Wolf, C.T. (2019b): "Professional identity and information use: on becoming a machine learning developer." *International Conference on Information*. Spring, Cham.

Wolf, C.T. (2020): "AI Models and Their Worlds: Investigating Data-Driven, AI/ML Ecosystems Through a Work Practices Lens." *International Conference on Information*. Springer, Cham.

van Someren, M. W., Y. F. Barnard, and J. A. C. Sandberg. (1994): "The think aloud method: a practical approach to modelling cognitive." *London: AcademicPress*.

van Wyk, K.R., and McGraw, G. (2005): Bridging the gap between software development and information security. *IEEE Security Privacy*. 3, 5 (Sep. 2005), 75–79. DOI:https://doi.org/10.1109/MSP.2005.118.

Viégas, F.B. and Wattenberg, M. (2008): TIMELINES Tag clouds and the case for vernacular visualization. *interactions*. 15, 4 (Jul. 2008), 49–52. DOI:https://doi.org/10.1145/1374489.1374501.

Zhang, Q., et al. (2018): "Unsupervised learning of neural networks to explain neural networks." *arXiv preprint arXiv:1805.07468*.