

What Did My AI Learn? How Data Scientists Make Sense of Model Behavior

ÁNGEL ALEXANDER CABRERA, Carnegie Mellon University, USA

MARCO TULLIO RIBEIRO, Microsoft Research, USA

BONGSHIN LEE, Microsoft Research, USA

ROBERT A. DELINE, Microsoft Research, USA

ADAM PERER, Carnegie Mellon University, USA

STEVEN M. DRUCKER, Microsoft Research, USA

Data scientists require rich mental models of how AI systems behave to effectively train, debug, and work with them. Despite the prevalence of AI analysis tools, there is no general theory describing how people make sense of what their models have learned. We frame this process as a form of sensemaking and derive a framework describing how data scientists develop mental models of AI behavior. To evaluate the framework, we show how existing AI analysis tools fit into this sensemaking process and use it to design AIFINITY, a system for analyzing image-and-text models. Lastly, we explored how data scientists use a tool developed with the framework through a think-aloud study with 10 data scientists tasked with using AIFINITY to pick an image captioning model. We found that AIFINITY's sensemaking workflow reflected participants' mental processes and enabled them to discover and validate diverse AI behaviors.

CCS Concepts: • **Human-centered computing** → **HCI theory, concepts and models**; *Visualization systems and tools*; • **Computing methodologies** → **Artificial intelligence**; *Computer vision*.

Additional Key Words and Phrases: machine learning, AI, machine behavior, machine learning testing, sense-making, visualization

ACM Reference Format:

Ángel Alexander Cabrera, Marco Tulio Ribeiro, Bongshin Lee, Robert A. DeLine, Adam Perer, and Steven M. Drucker. 2022. What Did My AI Learn? How Data Scientists Make Sense of Model Behavior. In *TOCHI '22: ACM Transactions on Computer-Human Interaction*. ACM, New York, NY, USA, Article 1, 27 pages. <https://doi.org/10.1145/3542921>

1 INTRODUCTION

Designers make sense of feedback to inform their designs [26], doctors make sense of health records to guide their diagnoses [84], and programmers make sense of code to debug their software [30]. Similarly, data scientists make sense of their machine learning (ML) or artificial intelligence (AI) models to improve their performance, decide when to use them, and analyze their real-world impacts. Having a thorough understanding of how an AI behaves is especially important to detect and mitigate serious concerns such as fairness [38] and safety [61] issues.

What does it mean to make sense of AI behavior? Let us explore the example of a data scientist who wants to make a website more accessible by including text descriptions (alt-text) for images. They find multiple AI services for captioning images and have to pick the option that works best for

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
TOCHI '22,

© 2022 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/10.1145/3542921>

their data. The data scientist compares the options by generating alt text with each AI for a sample of images and develops a mental model of how each AI *behaves*: which AI can describe certain activities, is better in low light, or is more grammatically accurate. With a deeper understanding of how each AI service behaves, the data scientist can decide which one to use for their data. This is just one example use case for understanding AI behavior, which is essential for tasks ranging from training new models to detecting dataset shift and mitigating real-world failures.

While important, behavioral analysis requires significant human attention to ideate, structure, and test hypotheses of AI behavior. Data scientists instead often resort to limited and ad hoc methods, such as manually testing edge cases or waiting for end-users to report failures of deployed models [1, 38, 41, 47]. A number of AI analysis tools aim to improve this process, including crowdsourcing methods for discovering failures [6, 13], algorithms for finding slices of data with high loss [21], and checklists of expected model behavior [68]. Although useful for specific tasks, these tools tend to only address portions of the analysis process and are hampered by challenges at other stages of the process. For example, methods for creating subgroups of data [14, 21, 42] do not tell the user *which* subgroups are the most important, while model checklists do not have mechanisms for discovering new behaviors.

This article introduces a sensemaking framework describing how data scientists develop mental models of AI behavior. By framing AI analysis as sensemaking, we aim to provide a language for describing AI analysis, help developers identify gaps in existing tooling, and encourage analysis tools supporting the full sensemaking process. Sensemaking is a well-established paradigm that describes how people structure the unknown by iteratively creating mental models from data [85]. To accurately describe AI analysis as sensemaking, we used abductive analysis to adapt Pirolli and Card [63]’s framework for data analysis to fit the steps specific to AI development gathered from empirical studies of practitioners. Our resulting framework (Figure 1) describes how people create mental models of AI behaviors by organizing instances into meaningful schemas and hypotheses. The *mental models* data scientists derive are their internal representations of the behaviors of a complex, often black-box, AI *model*.

We evaluated our framework across the three powers of interaction frameworks defined by Beaudouin-Lafon [9]: *descriptive*, *evaluative*, and *generative* power. To test the framework’s *descriptive* and *evaluative* power, how it can detail and compare a range of existing interfaces, we reviewed AI analysis tools and showed how they fit into the stages of our framework. We found that most tools only address half of the sensemaking process, either discovery tools for finding and organizing instances or evaluation tools for testing known behaviors. Systems that combine discovery and evaluation could help data scientists effectively validate newly discovered behaviors. Next, to directly test our framework’s *generative* power, the ability to inform new designs, we used it to create an AI analysis tool, AIFINITY, for exploring image-and-text models like visual question answering and image captioning. Image-and-text models have many complex behaviors, from stereotypes to grammar issues, that make them a challenging domain for AI analysis.

For our final evaluation of the framework, we explored how data scientists use a full sensemaking system. We conducted exploratory think-aloud studies with 10 professional data scientists tasked with using AIFINITY to choose between two image captioning AIs. Participants found that AIFINITY matched their mental process for understanding AI behavior, with some even independently describing their processes in sensemaking terms. Additionally, the complementary features helped participants find numerous significant behaviors and actively think about confirmation bias.

In summary, the main contributions of this work are the following.

- A **sensemaking framework** describing how people develop mental models of AI behavior.
- An **AI analysis tool** called AIFINITY designed using the framework.

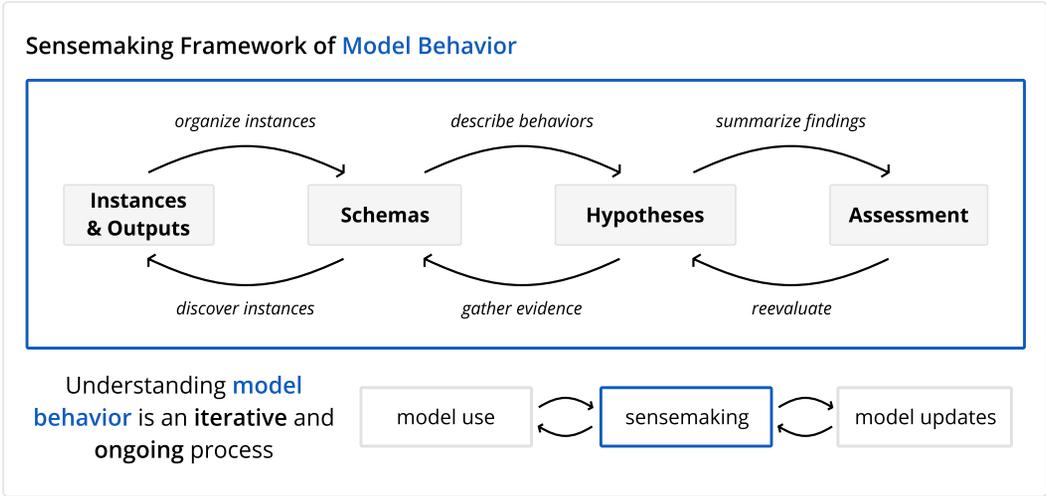


Fig. 1. The sensemaking framework describing how data scientists understand model behavior. We derived the framework from Pirolli and Card [63]’s sensemaking process and empirical studies of data scientist. The process is iterative and ongoing, with data scientists continuously reevaluating as they update and deploy their models.

- An **exploratory think-aloud study** with 10 professional data scientists to understand how people work with an AI analysis tool for the full sensemaking process.

2 BACKGROUND AND RELATED WORK

2.1 Behavioral Analysis of AI

Along with the growing use of AI systems in the real world, there has been increasing concern about the *behavior* of these systems [65]. Analyzing model behavior can uncover more nuanced, complex patterns not captured by aggregate metrics, such as how a model performs for particular subgroups or domains.

Behavioral analyses of AI systems in both academia and industry have discovered many real-world issues, some with significant societal implications. Fairness concerns are one major issue that plague many models trained on data about people. Notable biased systems include gender classification models that significantly underperformed for women of color [12] and criminal recidivism prediction models that classified people of color as higher risk [4]. Better understanding these models can also inform discussions of whether they should be used at all, such as models that use binary gender definitions [45]. Another area of concern is potential safety issues. For example, pedestrian detection systems may not work as well at night or in inclement weather [79], which when used in self-driving cars can lead to serious accidents [59]. Medical diagnosis models can have similarly serious errors, for example, a cancer screening model failing to detect a malignant tumor [61]. A growing number of researchers and practitioners are conducting deeper analyses of deployed AI systems to discover and mitigate these behaviors.

Knowing how an AI behaves can also be helpful in less critical settings. Improving an existing AI system often requires developers to know what types of data their model fails on so they can target their data collection [37]. Consistent failures can also indicate limitations of an AI model’s architecture and influence the design of future iterations.

Visual and algorithmic systems can help data scientists describe, detect, and validate the behavior of their models. These techniques range from tools for slicing and exploring model outputs to testing specific behaviors. We review several of these systems as we define the sensemaking framework in Section 4, and explore how they fit into different stages of the framework. By describing a theoretical framework of how data scientists understand model behaviors, we situate these existing systems in the broader analysis process and identify stages and domains with limited tooling.

2.2 Sensemaking

Sensemaking was originally formalized by Karl Weick, a social psychologist, in 1995 to describe how members of organizations come to a collective understanding of their surroundings [85]. At its most abstract, it can be thought of as “structuring the unknown,” or the “process through which individuals work to understand new, unexpected or confusing events” [55, 86]. It is an ongoing, iterative process by which people develop mental models of the world to make decisions and take actions. Weick’s formalization of sensemaking spurred numerous empirical and theoretical studies, ranging from how organizations work through crises [3] to how entrepreneurs deal with failure [81] and even how we should design explainable AI [44].

Sensemaking has since expanded beyond social psychology and has been applied to domains such as ecology [88] and medicine [17]. Most relevant to this work are the applications of sensemaking to HCI, where computer and information scientists framed data analysis as sensemaking: constructing a mental model from extensive unstructured data. One of the earliest formalizations came from Russell et al. [71], who defined a “learning-loop complex” in which analysts cycle between creating representations of a system and fitting data to those representations. Russell’s framework was later expanded by Pirolli and Card [63] to describe the specific steps and representations they observed data analysts use in practice.

Pirolli and Card [63]’s framework has become a frequent reference for data analysis and visualization research. One application of the framework has been structuring empirical studies of analysts, such as Grigoreanu et al. [30]’s study of programmers’ processes and challenges when debugging software. It has also been used to design data analysis tools, including visualizations for large graph networks [18] or tracking patterns in microblogs such as Twitter [11]. Researchers and developers have been able to create tools that better fit people’s processes by using Pirolli and Card [63]’s sensemaking framework.

In this work we adapt Pirolli and Card [63]’s sensemaking framework to AI analysis, as it is the most widely used framework in the most closely related domain, data analysis. As Pirolli and Card [63] did with Russell et al. [71]’s framework, we analyze empirical studies of AI practitioners to derive a new framework that more accurately describes the sensemaking process for understanding AI behavior. With a formal sensemaking framework specific to AI analysis, we hope to bring structure to the field, just as the above frameworks did in fields such as organizational psychology and data analysis.

3 METHODOLOGY

To create a framework that describes AI practitioners’ process we used *abductive analysis* [80] to iteratively adapt Pirolli and Card [63]’s sensemaking framework to empirical studies of AI/ML practitioners. In contrast to *inductive* methods such as grounded theory [78], which develop a framework from empirical evidence, and *deductive* approaches that directly apply existing theories, *abduction* extends or develops theory to explain new evidence. We decided that an abductive approach would be the most appropriate for this work since we adapt theory from a related domain, data analysis, to describe a new process, how practitioners understand AI behavior.

We primarily built from Pirolli and Card [63]’s sensemaking framework which describes how intelligence analysts make sense of large amounts of unstructured data. In their framing, analysts first go through an information foraging loop, where they filter *data sources* into a *shoebox* of relevant information. Snippets from documents in the shoebox make up the *evidence file*. Next is the core sensemaking loop, where analysts create *schemas*, structured organizations of the data, from the evidence file which are used to create and support *hypotheses*. Lastly, these hypotheses are used to create a final *presentation*. One can imagine a detective in front of a corkboard, cutting out and organizing newspaper clippings to pin them up and connect them with red thread.

To adapt Pirolli and Card [63]’s framework to AI analysis, we reviewed empirical studies of how practitioners work with AI systems in the real world. Since there are no survey papers, to date, directly covering this area, we relied primarily on academic search engines and citation graphs. Our review focused on studies with first-hand interviews and surveys to get the most direct look at data scientists’ processes (Table 1). For our analysis, we coded the empirical studies and used an affinity diagram to recursively fit the codes to the Pirolli and Card [63] sensemaking stages. During the abductive analysis, we also updated the stages to better describe AI practitioners’ processes. In the following section, we describe the resulting framework in detail and describe the key ways in which it differs from existing frameworks.

4 SENSEMAKING FRAMEWORK

The resulting sensemaking framework for understanding AI behavior is shown in Figure 1. The least structured stage is gathering **(1) instances and model outputs** from a variety of sources such as real-world users or synthetic methods. Data scientists then begin to organize the instances into general **(2) schemas** of semantically similar instances and behaviors. Schemas can be either rough groupings or strict slices of data. Data scientists then define formal **(3) hypotheses** of AI behaviors and gather additional evidence to validate their hypotheses. Lastly, data scientists derive a final **(4) assessment** of their discoveries, organizing hypotheses to be useful in subsequent tasks like choosing between AI services or updating a model’s architecture. The sensemaking process does

Table 1. Empirical studies of how practitioners work with AI systems. We synthesized insights from these studies to develop the sensemaking framework. We limited our search to papers that directly interviewed or surveyed AI/ML practitioners to study their real-world processes and challenges.

Study	Topic	Interview #	Survey #
Kim et al. [47]	Data Scientists in Software Teams	0	793
Wan et al. [82]	ML & Software Development	14	342
Serban et al. [73]	ML & Software Engineering	0	313
Holstein et al. [38]	ML Fairness in Industry	35	267
Zhang et al. [93]	Software Engineering & ML	8	195
Yang et al. [90]	Interactive ML	24	98
Sambasivan et al. [72]	Data Cascades in AI	53	0
Bhatt et al. [10]	XAI in Deployment	50	0
Hong et al. [40]	Human Factors & XAI	22	0
Muller et al. [57]	Data Scientists & Data	21	0
Hopkins and Booth [41]	AI Outside Big Tech	17	0
Nascimento et al. [58]	Development Processes in ML	7	0
Piorkowski et al. [62]	AI in Interdisciplinary Teams	4	0
	<i>total</i>	237	2,008

Table 2. How existing AI analysis systems fit into the sensemaking framework. Some of the tools focus on specific behaviors, like biases, or domains, like self-driving cars, but they all help data scientists better understand the behaviors of their AI systems at different points in the sensemaking process.

Venue	Paper	Instances	Schemas	Hypotheses	Assessment
<i>AAAI</i>	Beat the Machine [6]	■			
<i>arXiv</i>	Dynabench [46]	■			
<i>ICLR</i>	Goodfellow et al. [29]	■			
<i>CVPR</i>	StyleGAN [43]	■			
<i>JBD</i>	Data Augmentation [76]	■			
<i>VIS</i>	CAVA [16]	■			
<i>VLDB</i>	Snorkel [66]	■			
<i>WWW</i>	Patterned BTM [52]	■	■		
<i>VIS</i>	What-if Tool [87]	■	■		
<i>HCOMP</i>	Pandora [60]		■		
<i>AAAI</i>	Lakkaraju et al. [49]		■		
<i>arXiv</i>	Spotlight [23]		■		
<i>CVPR</i>	Barlow [77]		■		
<i>ICDE</i>	Slice Finder [21]		■		
<i>VIS</i>	FairVis [14]		■		
<i>CHI</i>	ModelTracker [2]		■		
<i>VIS</i>	Squares [67]		■		
<i>N/A</i>	Facets [64]		■		
<i>IUI</i>	AnchorVis [19]		■		
<i>HILDA</i>	MLCube [42]		■		
<i>CSCW</i>	Deblinder [13]	■	■	■	
<i>ACL</i>	Errudite [89]		■	■	
<i>ICLR</i>	Domino [24]		■	■	
<i>VIS</i>	HypoML [83]			■	
<i>ASE</i>	DeepRoad [92]			■	
<i>ICSE</i>	DeepTest [79]			■	
<i>ICSE</i>	Structure-Invariant Testing [35]			■	
<i>FAccT</i>	Interactive Model Cards [22]	■	■		■
<i>CHI</i>	Symphony [8]		■		■
<i>arXiv</i>	Robustness Gym [28]		■	■	■
<i>ACL</i>	Checklist [68]			■	■
<i>FAccT</i>	Model Cards [56]				■
<i>IBM JRD</i>	FactSheets [5]				■

not have to start from the initial stage of instances and outputs. Practitioners may have existing hypotheses, or may use tools that slice and organize instances into pre-defined schemas.

This adapted framework differs in a few key ways from the Pirolli and Card [63] formalization. Primarily, it is missing the initial foraging loop with the *shoebox* and *evidence file* stages. Unlike analysts who sort through data sources, such as newspapers, to extract snippets of evidence, AI analysis starts with instances, model inputs, that are directly relevant to a model’s behavior. While AI practitioners actively search for new instances to discover hypotheses, they do not have to further sort and modify instances in their sensemaking process. Next, the instances and outputs in AI analysis are lower level than the data sources, like research articles, used by analysts. Thus, the schemas for AI analysis tend to be groupings of instances rather than connections between high-level patterns or findings. This also means that hypotheses are directly verified using supporting instances and outputs, and need sufficient, diverse evidence to be accurately evaluated. Overall, the focus of AI analysis is on creating appropriate schemas and ensuring the validity of hypotheses rather than foraging for relevant evidence.

The context in which AI analysis occurs also differs significantly from sensemaking in domains such as data analysis. Sensemaking for AI systems is an iterative and ongoing process, as AI systems are constantly being updated and applied to new domains. In traditional data analysis, new reports or research may update existing hypotheses over time but often do not lead to brand new patterns. Updates to black-box AI systems, on the other hand, can completely change the behavior of an AI system and require reevaluating all hypotheses. Additionally, new instances are constantly being received from end users, informing new schemas and hypotheses. The volatility and quick iteration of AI systems have implications for tools that support the sensemaking process.

In the following sections, we describe in detail the four stages of the sensemaking process for AI analysis. In each section, we first describe how data scientists currently approach the sensemaking process and then describe existing tooling available at each stage.

4.1 Instances and Outputs

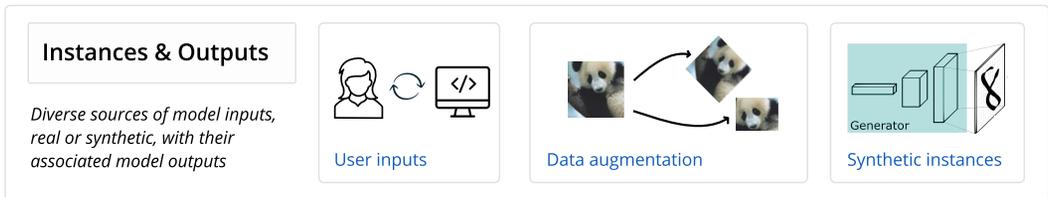


Fig. 2. The least structured stage of the sensemaking process consists of **instances and outputs**, model inputs from a variety of sources along with their associated model predictions. Instances can include both real-world user inputs and synthetic data.

At the core of the sensemaking process are data instances and their associated model predictions, the outputs of the model for the given instances (Figure 2). The most convenient source of instances are datasets collected to train an AI system, often split into training, validation, and testing sets on which aggregate metrics are calculated. While convenient, initial training datasets are limited and can lead to misleading performance measures and missed behaviors. For example, one participant in Wan et al. [82]’s study found significant overlap between their training and testing sets that produced an inflated model accuracy, while two participants interviewed by Hopkins and Booth [41] lamented that they needed a much greater diversity of instances than they had to accurately evaluate the performance of their model.

To better understand the behavior of their models, data scientists constantly collect new real-world instances to both update their models and discover new behaviors. This is especially important due to data drift, with 55% of the data scientists interviewed by Sambasivan et al. [72] describing factors such as new environmental factors and human patterns leading to model failures or unexpected outputs. The data scientists interviewed described monitoring the performance of the model over time on newly collected instances to identify performance drops or new regressions.

Despite the utility of real-world data, it is often expensive and slow to gather and label real instances, limiting developer access to data. Instead, data scientists “dogfood” their models, creating instances they think might be particularly difficult for an AI or show interesting behaviors [47]. Data scientists interviewed by Hopkins and Booth [41] found that this type of “prodding and probing” of models helped them better understand and work with black-box systems. Dogfood testing can be especially important for rare or sensitive behaviors which could have serious consequences in the real world [1].

Finally, it is not just the quantity and diversity of instances that is important for AI analysis, but what features are available for each instance. For example, to detect whether a model treats people of a certain demographic group inequitably, the data instances have to have a feature for that demographic information. Sensitive information, such as demographic details, is often not collected or present in a dataset and was one of the primary challenges for data scientists in discovering biases found by Holstein et al. [38]. In sum, both the number of instances *and* number of features of a dataset are important for discovering relevant behaviors.

4.1.1 Data collection and labeling methods. Tools at the instance and output stage often focus on scaffolding data collection, artificially generating instances, and adding features to a dataset.

Instead of waiting to gather real-world data from users, some techniques proactively use crowdworkers to gather instances. Beat the Machine (BTM) [6] and DynaBench [46] directly ask end-users to explicitly find instances for which a model fails, collecting instances that may surface interesting behaviors. Subsequent methods such as Deblinder Cabrera et al. [13] and Patterned Beat the Machine [52] build on this process by asking users to provide more context for a failure and find instances relevant for later schemas and hypotheses.

Data is often expensive to collect, so *synthetic*, artificially generated instances can provide a useful alternative to real-world instances. A common method for creating synthetic data is data augmentation, creating new instances by modifying existing ones, e.g., rotating or cropping images [76]. To create new instances that are not in a dataset, techniques like generative adversarial networks (GANs) can be used to generate novel examples [29]. StyleGAN is one such technique that generates new images from high-level semantic descriptions [43]. Synthetic instances are a low-cost way to augment a dataset, but it is not possible to generate any arbitrary instance, and synthetic instances are often less diverse than examples found in the real world.

There are also methods for adding new features to a dataset, providing details for each instance that can surface new behaviors. A separate AI model or heuristic functions are a common way to extract new features from an instance, such as the noisy labeling functions in Snorkel [66]. A related system is CAVA, which uses a knowledge graph to extract new attributes for an instance, such as populations from country names [16]. Additional features, or metadata, are essential for the subsequent stage of grouping and organizing instances into schemas.

Gathering diverse instances remains a challenging problem, as traditional methods remain expensive and synthetic techniques are noisy and limited to certain data types. In the context of the full sensemaking process, tools at the instance and output stage are often not informed by findings from later stages, such as interesting schemas or new hypotheses. For example, validating hypotheses requires collecting specific instances, which is often not well supported by current data

collection methods. Data collection or generation techniques that are more closely informed by the needs of schemas and hypotheses could better support data scientists' AI analysis process.

4.2 Schemas

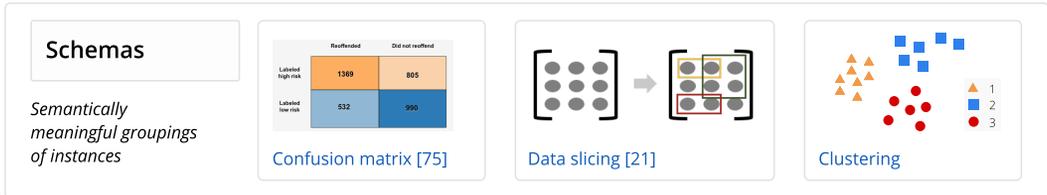


Fig. 3. Creating **schemas** is the second major sensemaking stage. Schemas are organizations of instances into meaningful layouts or groupings. Common schemas for AI outputs include confusion matrices, subgroups of data, and clusters.

The second sensemaking stage is organizing instances into semantically meaningful groups, called *schemas* [33, 63]. Schemas let practitioners hypothesize new model behaviors or collect evidence for existing hypotheses (Figure 3). There is significant flexibility in how schemas are created, from formal slices of a dataset to rough groupings of semantically similar instances.

Some of the most common schemas are classic methods for evaluating AI systems, such as the confusion matrix for classification problems [62, 75] and residual plots for regressions. Yang et al. [90] described the use of these visualizations as core knowledge required by the data scientists they spoke with. Splitting a model's output by predicted and ground truth output lets data scientists identify numerous metrics related to the model's behaviors; does the AI have a higher recall than precision? Is the false positive rate acceptable? These questions of model behavior are often central for data scientists, such as data scientists in Wan et al. [82]'s study making tradeoffs between metrics like precision and recall. Residual plots give a similar idea of how well a regression model behaves, as nonrandom errors can suggest a model is not adequately describing the data.

While these output-based visualizations may be helpful, they are limited to detecting behaviors described by output groups. Many important behaviors are found in groups defined by a model's *input* features; for example, fairness issues are defined by demographic information that is rarely the output of a model. Often called 'subgroup analysis,' or 'data slicing,' splitting and comparing instances by input features can detect such behaviors. Data scientists often look at model performance across these subgroups to track issues such as biases [38, 41].

For less structured data types such as images it can be challenging to create groups of similar instances in the first place, such as all images with a specific object in them. Without additional metadata collected or generated in the instances stage, it is not possible to create clear schemas for those semantic features. To address this, a data scientist in Holstein et al. [38]'s study wished for an oracle that would automatically find a hundred other examples of a failure they had found.

4.2.1 Creating schemas. There are myriad tools for creating and visualizing schemas of instances, from faceted layouts [64] to crowd-powered methods for finding areas of high error [60].

Better encodings of classic visualizations such as the confusion matrix can speed up and improve model analysis. For example, unit visualizations showing individual failures allow data scientists to dive deeper into the cause of low performance metrics [2, 67]. Confusion matrices can also be extended beyond binary classification, such as analyzing hierarchical models [31] or comparing multiple models [36].

Novel visualizations can be especially helpful for subgroup analysis. The most direct method is to look at groups of all combinations of features using, for example, data cube analysis [42]. Since this can create a countless number of subgroups, other visual systems allow users to create subgroups from specific features and values [14, 64, 87]. While useful if a data scientist knows what subgroups they want to create, these systems do not lead users towards interesting groupings. Automatic slicing algorithms such as Slice Finder can create a more reasonable number of subgroups with characteristics such as high loss [21]. By slicing data using input features, these visualizations and algorithms create schemas of subgroups highlighting important AI behaviors.

Beyond explicit data slicing, there are also tools for creating schemas of unstructured data. For example, clustering instances can surface semantically similar groups that may have interesting characteristics [7, 49]. Visualizations can also help semantically group data [89]; for example, AnchorVis [19] lets users define “anchors” that spread the data over different semantic dimensions.

Unfortunately, Holstein et al. [38] and Wan et al. [82] found that knowing *what* groups of instances to create and *how* to group instances are still major challenges for many data scientists. Current schema methods are mostly focused on highlighting known patterns in well-structured domains like tabular data. Additionally, few schema methods help data scientists move on to the hypothesis stage by formally defining hypotheses and gathering diverse supporting evidence. Schema methods that are better informed by hypotheses and can more meaningfully organize large, unstructured datasets could better support data scientists.

4.3 Hypotheses

Hypotheses

Formal descriptions of model behaviors with supporting evidence

Test case	Expected	Predicted	Pass?
A Testing Negation with MFT			
Labels: negative, positive, neutral			
Template: I (NEGATION) (POS VERB) che (THING).			
I can't say I recommend the food.	neg	pos	X
I didn't love the flight.	neg	neutral	X

[Checklists \[68\]](#)



[Test cases \[79\]](#)

Fig. 4. Creating **hypotheses** is the third sensemaking stage. Hypotheses are descriptions of model behavior with supporting evidence. Hypotheses can come from schemas or existing domain knowledge, like checklists and unit tests.

The third stage of the framework are hypotheses, formal descriptions of model behaviors (Figure 4). A hypothesis is a high-level description of a behavior (e.g., *the AI fails in low light*, or *the AI works best for long sentences*) along with supporting evidence. Data scientists test the validity of their hypotheses by gathering enough diverse data to determine how prevalent a behavior is. While hypotheses can come directly from schemas, they can also originate from a data scientist’s own domain knowledge or existing behaviors, such as a data scientist experienced with image models checking how a model performs in low-light settings.

Hypotheses in deployed settings are often described as unit or regression tests, well-defined tests of behavior hypotheses [38, 47, 93]. In some cases data scientists even use a test-driven ML approach in which they first define the behaviors that a model should have before training and evaluating the model [90]. For example, participants surveyed by Zhang et al. [93] often derive initial behaviors their models should have from specifications of the AI product they are developing. When updating their models, data scientists can check these hypotheses to ensure they are not regressing on important behaviors and monitor any improvements.

Varied external sources can provide hypotheses of model behavior, such as real-world users or customer service personnel. Looking through customer bug reports, customer-facing team members often go through the sensemaking process themselves, finding enough examples of an AI’s behavior to describe and report a hypothesis. Hong et al. [40] termed the people who find and test these hypotheses “model breakers”, roles who interact with customers and may have more direct knowledge of the ways in which a model may behave. From these initial hypotheses, data scientists or testing engineers can go back to the schema and instances stages to collect more evidence and validate the prevalence of reported hypotheses.

4.3.1 *Defining hypotheses.* Hypothesis tools help data scientists understand and test model behaviors, especially when tracking multiple hypotheses and assessing supporting evidence.

Visualization systems have shown promise for helping data scientists convert schemas into formal hypotheses. Errudite is a system for NLP models that lets data scientists slice their data into schemas *and* formally define hypotheses of model behavior [89]. Robustness Gym extends this capacity for NLP models by letting data scientists test a variety of hypotheses, from adversarial attacks to data augmentation [28]. There are also systems for statistical hypothesis testing, for example, HypoML is a visual system that lets data scientists statistically test how models perform across specific concepts [83].

Formal testing methods can help scaffold and evaluate hypotheses of model behavior. Even simple checklists of expected behaviors can give data scientists an idea of how well their AI performs in common scenarios [39, 68]. These checklists can be either general descriptions of behaviors or more specific hypotheses with supporting evidence that can validate if an AI shows a behavior. Similar to testing in software engineering, data scientists can also test more low-level behaviors of AI systems [91]. Metamorphic testing, checking if a permutation of an input has an expected impact on the output, can be used to test behaviors such as the impact of weather conditions on a self-driving car [92].

Current tools for creating and testing hypotheses tend to focus on specific, predefined behaviors. They often do not enable data scientists to go back to the schema and instances stages to discover new behaviors and hypotheses. There is also a more limited set of tools for this stage of the process compared to the schema stage. Robust hypothesis creation and evaluation tools could help data scientists more accurately describe and test what real-world behavior their models have.

4.4 Assessment

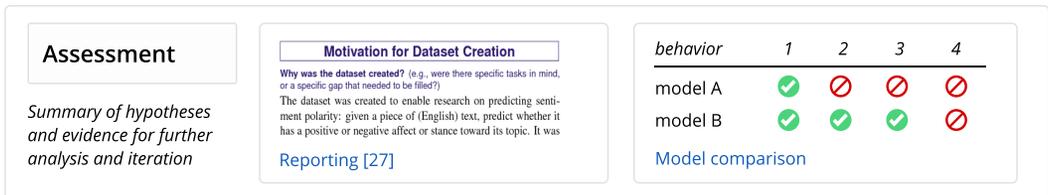


Fig. 5. An **assessment** of the model’s behavior is the final stage. The assessment provides an actionable summary of a model that can be used for tasks such as improving the model or choosing between different AI services.

Lastly, data scientists combine and organize hypotheses into a cohesive assessment of the behaviors of a model that can be used to make informed decisions (Figure 5). For example, when choosing between AI services, a data scientist needs a summary of the models’ behaviors to decide

which AI provides the best overall performance. Or, in AI development, ML practitioners need to know the most significant failures or areas in which their model can improve the most. Additionally, Yang et al. [90] found that ML consultants often report direct data insights and model iterations, assessments, to customers to increase their trust and reliance on a model.

As the most structured stage of the sensemaking process, assessments often act as the starting point for the other AI development processes. For example, a full assessment can be used to decide which AI service is the best for a certain dataset. It can also guide future data collection and model updates to target the areas for which the model performs the worst. Data scientists can then go back to the assessment to see how their updates have changed model behaviors.

AI teams often attempt to track model behaviors to check for serious issues and understand how their AI systems evolve over time. Many data science teams often deal with issues on a case-by-case basis, fixing problems as they are detected in the real world [38]. This introduces its own challenges of ensuring that model updates do not inadvertently regress on certain behaviors while improving others [82]. By having a combined central assessment of model behaviors, data scientists can quickly see their model's overall performance and make informed decisions [28, 68].

4.4.1 Assessment mediums. Recent work has explored how structured reporting about datasets and models can improve future iterations. For example, Datasheets for Datasets [56] tracks the metadata of a dataset, such as provenance and demographic distribution, to inform future model builders, while Model Cards [27] describe AI models to inform their use and potential downsides. Checklists of important steps and processes that data scientists should take can also lead data scientists to more proactively audit the behaviors of their models [54].

Most current assessment tools focus on aggregate metrics and characteristics of a model, whereas AI teams often end up tracking behaviors in an ad-hoc manner. Systems, especially visualizations, that can effectively summarize and track changes in behavior over time could provide a useful and actionable assessment for data scientists. This information can augment documentation methods, for example, with interactive model cards [22], and provide a holistic view of how an AI system is working. While assessment is the final sensemaking stage, it is not the end of the process. Understanding model behavior is an iterative and ongoing process that data scientists continue going through as they update their AI and see new behaviors in the real world.

5 AIFINFINITY SYSTEM

To assess our framework's *generative* power, we used it to create a system for analyzing image-and-text models called AIFINFINITY. AIFINFINITY can be used to understand the behavior of a single model using ground-truth labels or compare two models against each other. In the review of existing tools for AI analysis we found that there were a lack of systems that covered the full sensemaking process and helped data scientists move between sensemaking stages. Therefore, our aim was to design a system that met these two goals, using both new and existing AI analysis techniques.

We focused on image-and-text models since they are growing in use for tasks like image captioning, visual question answering, and optical character recognition. Although there are many tools for understanding the behavior of tabular and text models, as described in Section 4, there are few tools specifically for image models. Image data is often unstructured, making it difficult to explore instances and create meaningful schemas and hypotheses.

AIFINFINITY is a Jupyter widget written in Python and Typescript. Jupyter notebooks are one of the most common data science platforms for data analysis and model training [74]. By making AIFINFINITY a widget, we allow data scientists to directly load instances and model outputs from a computational notebook into the tool. AIFINFINITY is also model-agnostic, working with common AI platforms such as PyTorch, TensorFlow, and online services.

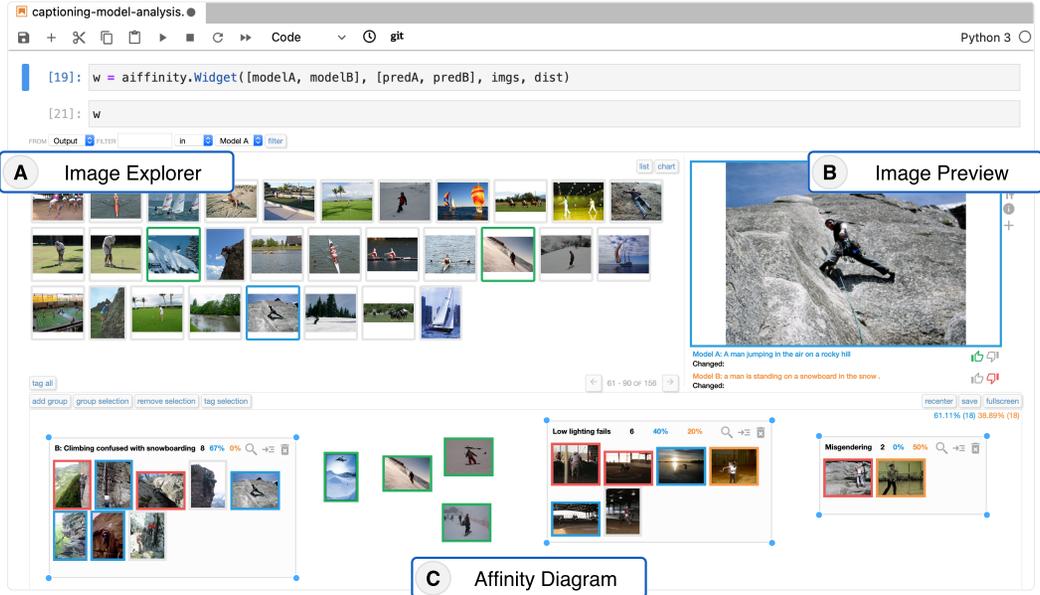


Fig. 6. The AIFINITY system is a Jupyter Widget that consists of three primary panels, shown here for the image captioning task used in the user study. The (A) **Image Explorer** shows a sample of images, sorted by the images with the most different labels. The (B) **Image Preview** shows the currently selected image. It lets users see the image’s extracted metadata, use tools find similar images, and create counterfactuals. Lastly, the (C) **Affinity Diagram** is where users can organize instances into schemas and hypotheses. The colored borders represent the quality judgements from users on whether they believe either or both of the outputs are adequate or not. Data scientists can load AIFINITY with any AI system and image dataset they are using in a Jupyter Notebook.

Running example: optical character recognition

AIFINITY supports various image-and-text models, but we focus on two primary examples for this work, optical character recognition (OCR) for the system walkthrough and image captioning for the user study in Section 6. As a running example of AIFINITY’s workflow, we walk through the example of an AI developer exploring whether their OCR system works for a new dataset of storefront signs [20]. This task is common in real-world scenarios such as Google Maps identifying the names of businesses from streetview data. As we describe AIFINITY, we use block quotes to describe how a data scientist could use each component in this running example (see Figure 7 for an overview):

Emma is an ML developer at a startup that provides an OCR service. Her company has a new client who wants to use the system to read street signs. Emma is unsure whether their model works for the client’s data, so she loads AIFINITY with a sample of the client’s storefront images, ground-truth labels, and the AI’s outputs. Her goal is to explore how well the AI works for this new dataset to decide whether she needs to collect new data and retrain the model.

5.1 Instances, Outputs, and Initial Schemas

AIFINITY is implemented as a Jupyter widget primarily to enable data scientists to use it with diverse, updating data, directly supporting the **instances and outputs** stage of the sensemaking process. Users pass to the widget a list or two of model outputs and image paths, which can be dynamically updated from the Jupyter notebook. AIFINITY explicitly supports two outputs for each instance for a couple of reasons. When analyzing a single model, one output can be the output of the AI model, while the other can be ground-truth labels. AIFINITY can also be used for model comparison, loading the outputs of both models. In both cases, comparing the two outputs provides a useful metadata feature for creating schemas and hypotheses.

The loaded images are displayed in AIFINITY's image explorer (Figure 6A), which shows them in a paginated list. When data scientists hover on a thumbnail or click to select an image, they see the full size version in the image preview (Figure 6B) on the right, along with the model output. AIFINITY initially sorts the instance exploration panel to show instances for which the two outputs are the most different. This creates an initial **schema** or grouping of the data that provides a sensible default for finding interesting hypotheses. When two outputs are significantly different, there is likely some interesting difference between the two. This technique is inspired by common loss functions for NLP models, namely the BLEU score for measuring sentence similarity [15], which we use to calculate how similar two labels are.

As data scientists discover interesting instances, they can drag them to the affinity diagram at the bottom of the interface (Figure 6C). Affinity diagrams are a common data analysis tool used in industry and research to organize and track data insights, especially in sensemaking processes [32]. Since images are two-dimensional and humans are especially good at 2D spatial cognition [50, 69], the affinity diagram is a compelling format for spatial organization of images. The affinity diagram serves two primary purposes in the AIFINITY system, allowing users to create rough **schemas** separate from the image list and to create and track **hypotheses** of behaviors.

As Emma explores the street sign images in the initial list, she finds that her model does not detect the text in a couple of round signs with text written in a circle. She drags these example images into the same area of the affinity diagram to keep track of them, creating an initial schema.

5.2 Schemas With Similar Search and Filtering

Beyond the initial sorted image list, AIFINITY provides a set of sorting and filtering tools to create new user-defined **schemas**. Since there is no direct technique to explore a dataset of images, unlike queries for tabular data, we provide two complementary features for creating new schemas, similar search and filtering.

AIFINITY's similar image search enables data scientists to discover instances that may have similar model behaviors. For a selected image in the image preview panel, a data scientist can click on the magnifying glass icon to find the most semantically similar images. Since pixels do not necessarily encode the semantic similarity of two images, AIFINITY instead uses the outputs of a pre-trained deep learning model to measure similarity. Specifically, AIFINITY runs each image through the ResNet-18 convolutional neural network (CNN) [34] trained on ImageNet and gets the second-to-last output layer, a 512-dimensional embedding vector representing the semantic content of the image. AIFINITY then calculates the cosine similarity between the selected image's embedding vector and all other images' vectors in the dataset and sorts the image exploration panel by the most similar images. The data scientist can then drag any interesting images into the affinity diagram. Similar image search acts as a **schema** of instances that are the most semantically similar to a reference image.

Emma wants to find more examples of round signs with text written in a circle. She selects the first image she found of a round sign and clicks the magnifying glass, which sorts the image explorer to show the most similar images. She finds various images of round signs that her model also fails to detect, so she drags them into the affinity diagram close to the original instance.

While similar image search is a useful heuristic for organizing instances, it is an approximate method that can be biased and miss related instances. AIFINFINITY lets data scientists filter images by various semantic features as a more formal way of schematizing the data. When the images are first loaded, AIFINFINITY runs two pre-trained deep learning models to extract metadata from the images. First, AIFINFINITY gets the ImageNet class of an image using the same pre-trained ResNet-18 model used for the similar image search. AIFINFINITY also runs an object detection model (FasterR-CNN ResNet-50 FPN [67]) trained on the MS-COCO dataset to extract common objects from the images. Data scientists can see the extracted metadata for a selected image by hovering over the information button to the right of the image in the image preview. In addition to the extracted metadata, they can also filter images by the labels of either source. For even more control, data scientists can also create custom tags for images that describe any feature of the image.

To filter images by any of these features, data scientists can use the *filter bar* at the top of the interface. Data scientists can use the filter bar to logically combine filters and isolate certain types of instances - for example, a data scientist could filter for images that have a certain object in them but do not have a keyword in the output. As data scientists add filters to the filter bar, the image exploration panel is updated to show only the matching images. Filtering is a **schema** that splits the dataset by explicit semantic features in contrast to similar search's rough grouping.

Emma hovers over the information button for a round sign with circular text and finds that it is incorrectly classified as an “analog clock”. While the class is incorrect, she thinks other round signs may have also been misclassified and decides to filter the images by the class “analog clock.” As expected, she finds various other round signs classified the same way, which she drags into the affinity diagram.

These image search and filtering techniques give data scientists multiple ways to schematize and mentally organize their data. From this general organization of images, they can then formulate and validate concrete hypotheses of AI behaviors.

5.3 Hypotheses and Assessment

In addition to being a medium for creating schemas of images, the affinity diagram also allows data scientists to create formal **hypotheses** of model behaviors. To create a hypothesis from the schemas, a data scientist can either select multiple images and click the “create hypothesis” button or drag the images into an existing hypothesis. Hypotheses are named rectangles in the affinity diagram data scientists can create for specific behaviors.

The initial evidence used to create a hypothesis is often not sufficient to fully support the prevalence of a behavior. To find more supporting evidence for a hypothesis, AIFINFINITY has a modified version of similar image search for hypotheses. When the magnifying glass on a hypothesis is clicked, AIFINFINITY calculates the average embedding vector of the images in the hypothesis and sorts the image exploration panel by the most similar images not already in the hypothesis. This allows data scientists to go back to the **schema** stage to find more supporting evidence.

To help data scientists get a more quantitative idea of how prevalent each behavior is, AIFINFINITY provides *quality judgements* that can be used to track whether an output is adequate for an instance. For each output on a given instance, a data scientist can indicate whether the output is correct

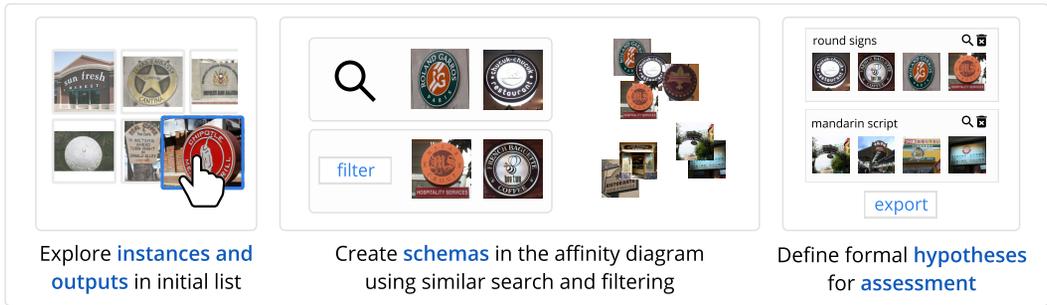


Fig. 7. An overview of the typical sensemaking process used by participants in the user study with AIFINITY. Participants often started by finding interesting instance in the Image Explorer. They then used the schema tools to find similar behaviors, and dragged them into the Affinity Diagram. Lastly, they created formal hypotheses from the schemas to find more evidence and organized their final assessment. The figure is shown for the optical character recognition task described in Section 5

by giving a thumbs up or thumbs down. Each hypothesis then shows the overall percentage of instances for which the data scientist indicated the labels are correct. This gives data scientists a quick quantitative view of how well their AI(s) perform for each hypothesis.

Emma has dragged various images of round signs into the affinity diagram and decides to create a formal hypothesis. She selects the images, clicks on the *create hypothesis* button, and names the resulting rectangle. To find more evidence, she uses the group similar image search by clicking the magnifying glass on the group. She finds a few more round signs and drags them into the hypothesis. She provides quality judgments for each image in the hypothesis and finds that her AI fails for more than 50% of the signs with circular text.

Since the original dataset may not have enough instances to adequately validate a hypothesis, AIFINITY also provides a counterfactual feature to allow the creation of more evidence and the refinement of hypotheses. Data scientists can click and drag to draw a black rectangle over an image in the image preview, occluding regions of the image to create a new instance. AIFINITY then runs the model on the newly modified image and shows the changed text output below the original output. The counterfactual tool allows data scientists to go back to the instances stage and create specific synthetic instances to test their hypotheses.

Most of the round signs with circular text that Emma found have logos in the center of the circle. Emma is worried that the AI system might actually be failing due to the logo, so she uses the counterfactual tool to create more evidence. She draws a black box in the center of a few of the images to remove the logos and adds the new images as evidence to her hypothesis. She finds that her AI is still not able to detect the text in the new images, further validating her hypothesis.

Lastly, data scientists can organize the affinity diagram with their evidenced conclusions into a final **assessment** of their model behavior, depending on the end goal of the analysis. These insights can then be saved and exported to share with other stakeholders and make actionable decisions.

Emma organizes the affinity diagram with the main hypotheses she has found, grouping them by the type and prevalence of each behavior. She exports the findings to save the results and uses them to improve her AI's performance for street signs by gathering more data and iterating on the AI's architecture.

6 USER STUDY

As a final evaluation of our framework, we conducted an exploratory think-aloud study with 10 professional data scientists tasked with using AIFINFINITY to choose between two image captioning models. This study aimed to understand how people use a complete sensemaking system, including how the different stages interact and how data scientists approach the process. We believe that these initial empirical insights can highlight the primary benefits and key features of AI analysis systems grounded in the sensemaking framework.

To recruit participants, we sent an email to 200 data scientists at Microsoft. We continued to invite participants in order of their responses until the qualitative themes in our iterative analysis converged at 10 participants (8 male, 2 female, mean age 32). The participants had an average of 6.8 years of data science experience and worked with various domains and models, including recommendation systems, search, captioning, and cybersecurity. The study lasted between 40 and 60 minutes, for which we compensated the participants with a \$25 Amazon gift card.

6.1 Study Procedure and Analysis

We started the study with a few background questions about the data scientist's experience with AI and behavioral analysis. The researcher then spent 10 to 20 minutes walking participants through AIFINFINITY, specifically for a task comparing two optical character recognition models used to read street signs, the same as the example in Section 5. The researcher explained the primary features and components of AIFINFINITY, and had the participant create at least one schema and hypothesis. We used a different domain and task for the introduction to not bias the behaviors that the participants looked for in the last part of the study.

In the final and main part of the study, which lasted 30 to 40 minutes, participants were tasked with using AIFINFINITY to choose between two image captioning models on a dataset of outdoor activities. This task was motivated by a common use case for image captioning, making photos accessible to people who are visually impaired or blind, for example, on social networks [53]. The task focused on model comparison to give participants a concrete goal, but since comparison requires participants to understand each model's behavior, our discoveries encompass understanding the behavior of one model. The first model, model A, was Microsoft's Cognitive Services image captioning system¹, and the second model, model B, was a pre-trained, off-the-shelf captioning model². Participants analyzed the behavior of the models on the UIUC Sports Event dataset [51], a collection of images from various indoor and outdoor sports. We chose this dataset as it has a wide variety of conditions, scenarios, and actions, while being a limited enough domain to explore in 30 to 40 minutes. To not limit or cherry pick the types of behaviors participants searched for, we gave them the general task of understanding the two models well enough to describe to a client, with supporting evidence, which model they should use for the given sports dataset.

As we conducted the studies, we transcribed the recordings and did iterative open coding of the results [70]. We also summarized the schemas and hypotheses of the participants as additional data on how the participants analyzed the two AI systems. With 10 participants, we found that the themes of how data scientists use a complete sensemaking system converged with significantly overlapping interaction patterns and hypotheses. After completing all the interviews, we conducted selective coding of the transcripts focused on the main themes identified in the open coding. We separate the findings into broader insights that are likely to generalize to other sensemaking systems and findings specific to the AIFINFINITY system.

¹<https://azure.microsoft.com/en-us/services/cognitive-services/computer-vision/>

²https://github.com/yunjey/pytorch-tutorial/tree/master/tutorials/03-advanced/image_captioning

6.2 Results

Making sense of model behavior

The challenges and goals described by the participants for AI analysis matched those identified in the empirical studies reviewed in Section 4. When describing their AI analysis workflow, all 10 participants talked about taking steps to better understand their AI systems beyond aggregate metrics. One participant (P8), a manager of an AI team, actually described their primary role as “*metric development*”: conducting behavioral analyzes on a deployed AI system and converting those insights into metrics to track and improve the system. Another participant (P5) described behavioral analysis as necessary because metrics like “*precision and recall can lie*”, but found that this deeper analysis is “*a very challenging problem.*”

Many of the strategies that participants use for AI analysis also reflect those described in the sensemaking framework. Five participants use human judges to label or gather instances, while two participants mostly rely on ad hoc spot checking like dogfooding to check if the AI is behaving as expected. Some data scientists have developed their systems for unit testing and validating model behaviors, with four participants using a form of “*regression sets*” that track specific model behaviors, or hypotheses. They use these sets to ensure that updates to their AI do not cause it to regress on important behaviors or subgroups of instances. Even the participants with bespoke tooling found behavioral analysis to be an open challenge, as one participant (P1) stated, “*we don’t really have a way of checking for patterns to see if a problem is a one-off or something more systematic.*” Like data scientists in the empirical studies, our participants tended to perform behavioral analysis in an ad hoc and post hoc manner, reacting to discovered failures.

Process and strategy

When the participants used the AIFINITY system, we noticed differences in how participants approached the sensemaking process. The first pattern we found was that participants started the AI analysis process from different stages. Since AIFINITY does not provide preexisting hypotheses, most of the participants (8) began their analysis by looking at the initial schema of instances with the largest output differences. The other two participants, who train image models in their work, started the analysis with their own preexisting hypotheses. They created these hypotheses from their experience and knowledge of how image models are most likely to fail. For example, a participant (P2) specifically created hypotheses for “*high contrast lighting*” and “*low light*” before looking at any of the instances. Despite starting at different stages, all participants eventually took an iterative process, going back to the image explorer to find new instances and using the affinity diagram to create schemas and hypotheses.

Another significant difference in participants’ processes was whether they took a breadth-first or depth-first approach. About half of the participants (4) took a breadth-first strategy by exploring multiple instances in the original schema before creating more specific schemas and hypotheses. The other six participants used a depth-first approach, immediately creating schemas and hypotheses for the first interesting instance they found. These different techniques led to a trade-off between the number of hypotheses and the amount of evidence participants found; participants using the breadth-first technique tended to find more hypotheses with less supporting evidence, while depth-first participants found fewer hypotheses with more evidence.

Complementary tools

One of the most salient benefits of having an integrated sensemaking system was the complementarity of tools across stages. As participants progressed through the sensemaking process, they had tools available to help them at each stage. For example, when participants wanted to validate an

initial idea of a behavior from a schema, they could create a hypothesis and find evidence using AIFINFINITY's similar image search feature. Participants found the progressions between tools and stages to be natural as they created schemas and validated hypotheses.

An unexpected benefit of AIFINFINITY was the complementarity of the features *within* each sense-making stage. This complementarity was most apparent in the schema stage with similar search and filtering tools. The benefit of having both tools was highlighted by one participant (P9), who in validating the hypothesis that models could not describe large groups of people found that *“using the tool together is useful, because otherwise, I was trying to look at [images with] groups of people but [similar image search] didn't give me that, but the object detection model is more specific.”* Similar search is a less structured but quicker schema tool, while filtering can create more specific and structured schemas. Participants generally started with the similar search tool to get an initial group of instances for a schema but were concerned about missing evidence with the “black box” search and so moved on to use the filtering approach. Having a quick heuristic tool combined with a more deliberate schema method was an essential feature of AIFINFINITY.

Dealing with confirmation bias

Confirmation bias is a significant challenge when creating and validating any hypothesis; How does a data scientist know that they have enough diverse instances to support their hypothesis? We found that having a combined sensemaking system helped data scientists combat confirmation bias. This was especially true when participants went from the hypothesis stage back to the schema stage to find more evidence, as they had various techniques at their disposal to discover or create more evidence. Six of the 10 participants found that at least one of their hypotheses did not hold after finding additional evidence. For example, a participant (P8) thought model A typically confused racquets for video game controllers, but quickly disproved their hypothesis by using the similar image search to find more images of people with racquets that were correctly described. Three participants also actively reflected on their potential confirmation bias and took steps to counteract it by proactively looking for disconfirming evidence.

Actionable, evidenced hypotheses

Overall, the participants found various hypotheses with significant supporting evidence. Participants created 4.1 hypotheses on average, which ranged from specific failures to high-level patterns. The most specific hypotheses included *“model cannot describe images with cliff backgrounds,”* and *“model fails to describe large groups of people on boats.”* Some of the most general hypotheses included *“model doesn't describe the central activity,”* *“the model is often too vague,”* and *“bad lighting leads to inaccurate captions.”* There was significant overlap in the hypotheses and behaviors the participants discovered despite the wide range of described behaviors. Five of the 10 participants found that Model B confused climbing images with snow, skiing or snowboarding. Four participants found that both models described most of the racquet sports as tennis and did not have *badminton* in their language. Lastly, the most common groupings eight participants created were for a specific activity, for example *climbing, boats, or tennis.*

At the end of the study, most participants had developed nuanced conclusions about which model they would choose for a given task. The most common conclusion, which seven of the participants came to, was that model A is more conservative, less detailed, but often correct, while model B provides more detailed captions, but is often wrong. Given these findings, they decided to make different recommendations for which model should be used depending on the risk profile and domain of the client.

Beyond describing the differences between the two models, some participants also asked questions about the underlying model and data and came up with potential fixes for the issues they

saw. Three participants attributed the patterns they found to biases in the training data or labels. Two of these participants hypothesized that there might be an “alpine” or “snow” bias in the data, causing model B to describe people climbing as snowboarding or skiing, and they wanted to look at the training data to verify their hypotheses. Two other participants hypothesized that the models themselves may be causing the problem by not having certain words in their vocabulary, specifically “badminton” and “croquet”, which were often described as “tennis” and “baseball.”

Using the AIFINFINITY system

We also found insights specific to the AIFINFINITY system and analysis of image and text models. Participants generally found the affinity diagram to be intuitive and usable, with five participants specifically stating that it was their favorite part of the interface and one participant (P3) stating that it “*makes total sense, especially for images.*” One participant (P8) especially liked the split between the top and bottom areas of AIFINFINITY, seeing them as two different representations of the data, or schemas: “*Switching between text and visual representations is very interesting - I can have a hypothesis and go back and forth.*” Affinity diagramming is a prolific sensemaking tool in other domains [32], which lends another piece of support to taking a sensemaking lens to AI analysis.

A feature that received mixed feedback in AIFINFINITY was the thumbs up or down quality judgement. Two participants (P1, P9) used it as the primary way of tracking which model was performing better, and a third participant (P2) liked that “*having them [images] colored gives you a quantitative feel for how strong your hypothesis is or not.*” While more than half of the participants (6) liked to have a quantitative view of their findings, four participants found that the judgement was too coarse to be very useful. Both captions were often wrong, but one was slightly better, or a caption being ‘good’ would depend on the situation. The participants would have liked more detailed descriptions to capture these nuances, such as scale or text descriptions.

Participants thought the counterfactual feature was useful but found that AIFINFINITY’s implementation of drawing black boxes was too simple. The participants wanted more image manipulation tools, such as adding new objects or changing image properties. Counterfactuals are a powerful tool for generating more evidence, and participants wanted these improved interactions to test more nuanced and complex behaviors.

7 DISCUSSION AND FUTURE WORK

Through our review of existing studies and tools, the design of AIFINFINITY, and the exploratory user study, we found that the sensemaking lens adequately describes how data scientists analyze AI systems. By describing AI analysis using a formal framework, we hope to give researchers and tool creators the language to better understand the context of their systems and studies in data scientists’ overall process. Future work can aim to fill tooling gaps for AI analysis or better understand the challenges and trade-offs in the different sensemaking stages.

7.1 Applications and Extensions of AIFINFINITY

The think-aloud study focused on one application of AIFINFINITY, comparing AI services, but AIFINFINITY can be used in various real-world AI analysis scenarios. When working with one model, data scientists can use AIFINFINITY to supplement traditional evaluation methods, such as aggregate metrics, by discovering, formalizing, and testing specific model behaviors on a labeled dataset. An example of this process is described in Section 5, in which a data scientist tests their model on a new dataset. When using AIFINFINITY for model comparison, data scientists can use it on their models, comparing iterations of an AI system using a new architectures or training set.

Participants found the initial set of schemas and hypothesis testing tools to be useful, but additional tools would have to be added to AIFINFINITY to make it widely applicable to real-world

models. Specifically, there were various behaviors that participants were unable to validate with the current tool set. This was especially the case for issues regarding bias and fairness, which participants wanted to test for but AIFINFINITY does not explicitly support. For example, creating schemas across demographic information can help identify potential biases, but the existing metadata did not have those features.

AIFINFINITY is primarily an exploratory analysis tool for image models and relatively small datasets that may not generalize to other use cases and domains. The affinity diagramming-based interface can work for other visual data such as videos but may not be adequate for encoding other data types such as text or audio. AIFINFINITY also requires users to manually select, explore, and organize instances, which cannot be manually done on a scale of thousands or millions of instances. For formal hypothesis testing on large datasets, especially when comparing models over time, a different system or extensions to AIFINFINITY would be required.

7.2 Gaps in Current Tooling

In reviewing the current landscape of AI analysis systems, we found a few significant gaps in current tooling. The first limitation is the lack of connection between the “discovery” half (*instances* ↔ *schemas*) and the “evaluation” half (*hypotheses* ↔ *assessment*) of the sensemaking process. There are many systems focused on the discovery half of the process that help data scientists slice and explore their data, and many systems for the evaluation half, like checklists and unit tests, that let data scientists validate known behaviors. There are comparatively few tools that let data scientists move between these two processes - turning rough groupings into formal hypotheses or discovering new hypotheses to validate.

There is also a lack of tools for certain types of data domains. Most schema and hypotheses tools are designed for tabular, timeseries, or text data that can be easily sliced and grouped. Unstructured data, such as images and videos, are harder to organize and there exist few usable tools for those domains in most stages of the sensemaking process. With the growing use of image and video recognition in the real world, behavioral analysis tools will be important in detecting and describing their behavior, especially for potential safety or fairness concerns.

7.3 Designing and Evaluating Tools With the Framework

The initial patterns found in the user study have some implications for future system design and empirical studies. For example, we found that there is a trade-off between using a breadth vs. depth-first approach when analyzing AI behaviors. A breadth-first approach tends to generate more hypotheses with less evidence, while a depth-first approach leads to fewer hypotheses with more supporting instances. Further experiments or studies could explore whether this leads to disparate insights and whether or not analysis systems should guide data scientists toward balancing these strategies. Other differences in approaches, like starting from certain sensemaking stages, could also be studied to improve data scientist processes.

The process of AI analysis also interacts significantly with other AI tools and processes. For example, explainable AI can be a useful tool at different points in the sensemaking process to both discover and evaluate hypotheses. Model updates and iteration also directly interplay with sensemaking, as people have to make sense of an updated model’s new or changed behaviors. These are both complex fields and topics which we did not explore in depth but which interact significantly with understanding AI behavior. Further studies of data scientists and deeper explorations of these interactions could identify areas where tools could bridge or better connect processes, for example, quick feedback loops between model updates and behavioral analysis.

Sensemaking has been applied to domains ranging from organizational psychology to data analysis. Each of these fields has developed unique techniques and tools to improve sensemaking

processes that could be used as inspiration for improving AI analysis. For example, there is a growing body of work on *distributed* or *crowd* sensemaking [25, 48], aggregating and reusing schemas and hypotheses from multiple people. Future work could explore how these concepts could be applied to improve AI analysis, for example, reusing common schemas and hypotheses between datasets and models.

7.4 Limitations

It is challenging to validate the usefulness of a theoretical framework, and our initial evaluation inherently has some limitations. First, when reviewing existing studies and systems, we likely missed some work that covers stages of our framework or fits the sensemaking process. While we do not claim that we conducted an exhaustive review of the literature, we believe that we covered the major works and subfields of AI analysis rigorously enough to support our framework. Second, to test the *generative* power of our framework, we implemented only one system for the specific domain of image and text models. Although it was not feasible to create multiple sensemaking systems, we believe that the reviewed systems provide a strong foundation for the framework, while the implementation of AIFINFINITY serves as a case study of how a complete sensemaking tool can be created. Lastly, our think-aloud study was conducted with participants at one company. While some of their procedures and the insights we derived may have been specific to that company's processes, we chose participants from different teams and suborganizations that act independently in order to increase the generalizability of our results.

8 CONCLUSION

This work introduces a sensemaking framework that describes how practitioners develop mental models of AI behavior. We derived the framework using a sensemaking lens and empirical studies of AI/ML practitioners. We then designed and implemented AIFINFINITY, an interactive tool for analyzing image-and-text models, and explored the dynamics of the sensemaking process in an exploratory think-aloud study with 10 professional data scientists. Researchers, designers, and tool creators can use the framework to better understand how people analyze AI systems and develop systems that are grounded in data scientists' analysis process.

ACKNOWLEDGMENTS

We thank Tongshuang Wu, Fred Hohman, and Will Epperson for their support and insight. We also thank the data scientists at Microsoft who participated in our studies. This material is based on work supported by the National Science Foundation Graduate Research Fellowship Program under Grant No. DGE1745016. Any opinions, findings, conclusions, or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

REFERENCES

- [1] Saleema Amershi, Andrew Begel, Christian Bird, Robert DeLine, Harald Gall, Ece Kamar, Nachiappan Nagappan, Besmira Nushi, and Thomas Zimmermann. 2019. Software Engineering for Machine Learning: A Case Study. In *2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*. IEEE, Montreal, QC, Canada, 291–300. <https://doi.org/10.1109/ICSE-SEIP.2019.00042>
- [2] Saleema Amershi, Max Chickering, Steven M. Drucker, Bongshin Lee, Patrice Simard, and Jina Suh. 2015. ModelTracker: Redesigning Performance Analysis Tools for Machine Learning. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*. ACM, Seoul Republic of Korea, 337–346. <https://doi.org/10.1145/2702123.2702509>
- [3] Deborah Ancona. 2012. Sensemaking: Framing and Acting in the Unknown. *The Handbook for Teaching Leadership: Knowing, doing, and being* (2012), 3–19. <https://doi.org/10.5465/amle.2011.0007>
- [4] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2016. Machine Bias — ProPublica.

- [5] M. Arnold, D. Piorowski, D. Reimer, J. Richards, J. Tsay, K. R. Varshney, R. K. E. Bellamy, M. Hind, S. Houde, S. Mehta, A. Mojsilovic, R. Nair, K. Natesan Ramamurthy, and A. Olteanu. 2019. FactSheets: Increasing trust in AI services through supplier’s declarations of conformity. *IBM Journal of Research and Development* 63, 4/5 (July 2019), 6:1–6:13. <https://doi.org/10.1147/JRD.2019.2942288>
- [6] Josh Attenberg, Panagiotis G. Ipeirotis, and Foster Provost. 2011. Beat the Machine: Challenging Workers to Find the Unknown Unknowns.
- [7] Gagan Bansal and Daniel S. Weld. 2018. A Coverage-Based Utility Model for Identifying Unknown Unknowns. In *32nd AAAI Conference on Artificial Intelligence, AAAI 2018*. 1463–1470.
- [8] Alex Bäuerle, Ángel Alexander Cabrera, Fred Hohman, Megan Maher, David Koski, Xavier Suau, Titus Barik, and Dominik Moritz. 2022. Symphony: Composing Interactive Interfaces for Machine Learning. In *CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (*CHI '22*). Association for Computing Machinery, New York, NY, USA, Article 210, 14 pages. <https://doi.org/10.1145/3491102.3502102>
- [9] Michel Beaudouin-Lafon. 2004. Designing Interaction, not Interfaces. In *Proceedings of the working conference on Advanced visual interfaces - AVI '04*. ACM Press, Gallipoli, Italy, 15. <https://doi.org/10.1145/989863.989865>
- [10] Umang Bhatt, Alice Xiang, Shubham Sharma, Adrian Weller, Ankur Taly, Yunhan Jia, Joydeep Ghosh, Ruchir Puri, José M. F. Moura, and Peter Eckersley. 2020. Explainable Machine Learning in Deployment. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. ACM, Barcelona Spain, 648–657. <https://doi.org/10.1145/3351095.3375624>
- [11] Harald Bosch, Dennis Thom, Florian Heimerl, Edwin Puttmann, Steffen Koch, Robert Kruger, Michael Worner, and Thomas Ertl. 2013. ScatterBlogs2: Real-Time Monitoring of Microblog Messages through User-Guided Filtering. *IEEE Transactions on Visualization and Computer Graphics* 19, 12 (Dec. 2013), 2022–2031. <https://doi.org/10.1109/TVCG.2013.186>
- [12] Joy Buolamwini and Timnit Gebru. 2018. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency (Proceedings of Machine Learning Research, Vol. 81)*, Sorelle A. Friedler and Christo Wilson (Eds.). PMLR, 77–91. <https://proceedings.mlr.press/v81/buolamwini18a.html> Buolamwini2018.
- [13] Ángel Alexander Cabrera, Abraham J. Druck, Jason I. Hong, and Adam Perer. 2021. Discovering and Validating AI Errors With Crowdsourced Failure Reports. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW2, Article 425 (oct 2021), 22 pages. <https://doi.org/10.1145/3479569>
- [14] Ángel Alexander Cabrera, Will Epperson, Fred Hohman, Minsuk Kahng, Jamie Morgenstern, and Duen Horng Chau. 2019. FairVis: Visual Analytics for Discovering Intersectional Bias in Machine Learning. In *2019 IEEE Conference on Visual Analytics Science and Technology (VAST)*. 46–56. <https://doi.org/10.1109/VAST47406.2019.8986948>
- [15] Chris Callison-Burch, Miles Osborne, and Philipp Koehn. 2006. Re-evaluating the Role of Bleu in Machine Translation Research. In *11th Conference of the European Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics, Trento, Italy, 249–256. <https://aclanthology.org/E06-1032>
- [16] Dylan Cashman, Shenyu Xu, Subhajit Das, Florian Heimerl, Cong Liu, Shah Rukh Humayoun, Michael Gleicher, Alex Endert, and Remco Chang. 2021. CAVA: A Visual Analytics System for Exploratory Columnar Data Augmentation Using Knowledge Graphs. *IEEE Transactions on Visualization and Computer Graphics* 27, 2 (Feb. 2021), 1731–1741. <https://doi.org/10.1109/TVCG.2020.3030443>
- [17] Joseph Chee Chang, Nathan Hahn, and Aniket Kittur. 2020. Mesh: Scaffolding Comparison Tables for Online Decision Making. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology*. ACM, Virtual Event USA, 391–405. <https://doi.org/10.1145/3379337.3415865>
- [18] Duen Horng Chau, Aniket Kittur, Jason I. Hong, and Christos Faloutsos. 2011. Apollo: Making Sense of Large Network Data by Combining Rich User Interaction and Machine Learning. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, Vancouver BC Canada, 167–176. <https://doi.org/10.1145/1978942.1978967>
- [19] Nan-Chen Chen, Jina Suh, Johan Verwey, Gonzalo Ramos, Steven Drucker, and Patrice Simard. 2018. AnchorViz: Facilitating Classifier Error Discovery through Interactive Semantic Data Exploration. In *23rd International Conference on Intelligent User Interfaces*. ACM, Tokyo Japan, 269–280. <https://doi.org/10.1145/3172944.3172950>
- [20] Chee Kheng Ch’ng and Chee Seng Chan. 2017. Total-Text: A Comprehensive Dataset for Scene Text Detection and Recognition. In *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, Kyoto, 935–942. <https://doi.org/10.1109/ICDAR.2017.157>
- [21] Yeounoh Chung, Tim Kraska, Neoklis Polyzotis, Ki Hyun Tae, and Steven Euijong Whang. 2019. Slice Finder: Automated Data Slicing for Model Validation. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*. IEEE, Macao, Macao, 1550–1553. <https://doi.org/10.1109/ICDE.2019.00139>
- [22] Anamaria Crisan, Margaret Drouhard, Jesse Vig, and Nazneen Rajani. 2022. Interactive Model Cards: A Human-Centered Approach to Model Documentation. *arXiv:2205.02894 [cs]* (May 2022). <https://doi.org/10.1145/3531146.3533108> arXiv: 2205.02894.

- [23] Greg d'Eon, Jason d'Eon, James R. Wright, and Kevin Leyton-Brown. 2021. The Spotlight: A General Method for Discovering Systematic Errors in Deep Learning Models. *arXiv:2107.00758 [cs, stat]* (Oct. 2021). <http://arxiv.org/abs/2107.00758> arXiv: 2107.00758.
- [24] Sabri Eyuboglu, Maya Varma, Khaled Saab, Jean-Benoit Delbrouck, Christopher Lee-Messer, Jared Dunmon, James Zou, and Christopher Ré. 2022. Domino: Discovering Systematic Errors with Cross-Modal Embeddings. *arXiv:2203.14960 [cs]* (April 2022). <http://arxiv.org/abs/2203.14960> arXiv: 2203.14960.
- [25] Kristie Fisher, Scott Counts, and Aniket Kittur. 2012. Distributed Sensemaking: Improving Sensemaking by Leveraging the Efforts of Previous Users. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, Austin Texas USA, 247–256. <https://doi.org/10.1145/2207676.2207711>
- [26] Eureka Foong, Darren Gergle, and Elizabeth M. Gerber. 2017. Novice and Expert Sensemaking of Crowdsourced Design Feedback. *Proceedings of the ACM on Human-Computer Interaction* 1, CSCW (Dec. 2017), 1–18. <https://doi.org/10.1145/3134680>
- [27] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé, and Kate Crawford. 2018. Datasheets for Datasets. (2018). <http://arxiv.org/abs/1803.09010>
- [28] Karan Goel, Nazneen Rajani, Jesse Vig, Samson Tan, Jason Wu, Stephan Zheng, Caiming Xiong, Mohit Bansal, and Christopher Ré. 2021. Robustness Gym: Unifying the NLP Evaluation Landscape. (2021), 1–34. <http://arxiv.org/abs/2101.04840>
- [29] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and harnessing adversarial examples. *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings* (2015), 1–11.
- [30] Valentina Grigoreanu, Margaret Burnett, Susan Wiedenbeck, Jill Cao, Kyle Rector, and Irwin Kwan. 2012. End-user Debugging Strategies: A Sensemaking Perspective. *ACM Transactions on Computer-Human Interaction* 19, 1 (March 2012), 1–28. <https://doi.org/10.1145/2147783.2147788>
- [31] Jochen Görtler, Fred Hohman, Dominik Moritz, Kanit Wongsuphasawat, Donghao Ren, Rahul Nair, Marc Kirchner, and Kayur Patel. 2022. Neo: Generalizing Confusion Matrix Visualization to Hierarchical and Multi-Output Labels. In *CHI Conference on Human Factors in Computing Systems*. ACM, New Orleans LA USA, 1–13. <https://doi.org/10.1145/3491102.3501823>
- [32] Gunnar Harboe and Elaine M. Huang. 2015. Real-World Affinity Diagramming Practices: Bridging the Paper-Digital Gap. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*. ACM, Seoul Republic of Korea, 95–104. <https://doi.org/10.1145/2702123.2702561>
- [33] Stanley G Harris and Stanley G Harris. 1994. Organizational Culture and Individual Sensemaking : A Schema-based Perspective. *INFORMS* 5, 3 (1994), 309–321.
- [34] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Las Vegas, NV, USA, 770–778. <https://doi.org/10.1109/CVPR.2016.90>
- [35] Pinjia He, Clara Meister, and Zhendong Su. 2020. Structure-Invariant Testing for Machine Translation. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. ACM, Seoul South Korea, 961–973. <https://doi.org/10.1145/3377811.3380339>
- [36] Andreas Hinterreiter, Peter Ruch, Holger Stitz, Martin Ennemoser, Jurgen Bernard, Hendrik Strobel, and Marc Streit. 2020. ConfusionFlow: A model-agnostic visualization for temporal analysis of classifier confusion. *IEEE Transactions on Visualization and Computer Graphics* (2020), 1–1. <https://doi.org/10.1109/TVCG.2020.3012063>
- [37] Fred Hohman, Kanit Wongsuphasawat, Mary Beth Kery, and Kayur Patel. 2020. Understanding and Visualizing Data Iteration in Machine Learning. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. ACM, Honolulu HI USA, 1–13. <https://doi.org/10.1145/3313831.3376177>
- [38] Kenneth Holstein, Jennifer Wortman Vaughan, Hal Daumé, Miro Dudik, and Hanna Wallach. 2019. Improving Fairness in Machine Learning Systems: What Do Industry Practitioners Need?. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM, Glasgow Scotland Uk, 1–16. <https://doi.org/10.1145/3290605.3300830>
- [39] Matthew K. Hong, Adam Fournay, Derek DeBellis, and Saleema Amershi. 2021. Planning for Natural Language Failures with the AI Playbook. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. ACM, Yokohama Japan, 1–11. <https://doi.org/10.1145/3411764.3445735>
- [40] Sungsoo Ray Hong, Jessica Hullman, and Enrico Bertini. 2020. Human Factors in Model Interpretability: Industry Practices, Challenges, and Needs. *Proceedings of the ACM on Human-Computer Interaction* 4, CSCW1 (May 2020), 1–26. <https://doi.org/10.1145/3392878>
- [41] Aspen Hopkins and Serena Booth. 2021. Machine Learning Practices Outside Big Tech: How Resource Constraints Challenge Responsible Development. In *AIES 2021 - Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. Association for Computing Machinery, Inc, 134–145. <https://doi.org/10.1145/3461702.3462527>
- [42] Minsuk Kahng, Dezhi Fang, and Duen Horng (Polo) Chau. 2016. Visual Exploration of Machine Learning Results Using Data Cube Analysis. In *Proceedings of the Workshop on Human-In-the-Loop Data Analytics - HILDA '16*. ACM Press, San

- Francisco, California, 1–6. <https://doi.org/10.1145/2939502.2939503>
- [43] Tero Karras, Samuli Laine, and Timo Aila. 2019. A Style-Based Generator Architecture for Generative Adversarial Networks. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Long Beach, CA, USA, 4396–4405. <https://doi.org/10.1109/CVPR.2019.00453>
- [44] Harmanpreet Kaur, Eytan Adar, Eric Gilbert, and Cliff Lampe. 2022. Sensible AI: Re-imagining Interpretability and Explainability using Sensemaking Theory. *arXiv:2205.05057 [cs]* (May 2022). <https://doi.org/10.1145/3531146.3533135> arXiv: 2205.05057.
- [45] Os Keyes. 2018. The Misgendering Machines: Trans/HCI Implications of Automatic Gender Recognition. *Proceedings of the ACM on Human-Computer Interaction* 2, CSCW (Nov. 2018), 1–22. <https://doi.org/10.1145/3274357>
- [46] Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Adina Williams. 2021. Dynabench: Rethinking Benchmarking in NLP. *arXiv:2104.14337 [cs]* (April 2021). <http://arxiv.org/abs/2104.14337> arXiv: 2104.14337.
- [47] Miryung Kim, Thomas Zimmermann, Robert DeLine, and Andrew Begel. 2018. Data Scientists in Software Teams: State of the Art and Challenges. *IEEE Transactions on Software Engineering* 44, 11 (Nov. 2018), 1024–1038. <https://doi.org/10.1109/TSE.2017.2754374>
- [48] Aniket Kittur, Andrew M. Peters, Abdigani Diriye, and Michael Bove. 2014. Standing on the schemas of giants: socially augmented information foraging. In *Proceedings of the 17th ACM conference on Computer supported cooperative work & social computing*. ACM, Baltimore Maryland USA, 999–1010. <https://doi.org/10.1145/2531602.2531644>
- [49] Himabindu Lakkaraju, Ece Kamar, Rich Caruana, and Eric Horvitz. 2017. Identifying Unknown Unknowns in the Open World: Representations and Policies for Guided Exploration. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence (AAAI'17)*. AAAI Press, 2124–2132.
- [50] Fritz Lekschas, Xinyi Zhou, Wei Chen, Nils Gehlenborg, Benjamin Bach, and Hanspeter Pfister. 2021. A Generic Framework and Library for Exploration of Small Multiples through Interactive Piling. *IEEE Transactions on Visualization and Computer Graphics* 27, 2 (Feb. 2021), 358–368. <https://doi.org/10.1109/TVCG.2020.3028948>
- [51] Tianyi Li, Kurt Luther, and Chris North. 2018. CrowdIA: Solving Mysteries with Crowdsourced Sensemaking. *Proceedings of the ACM on Human-Computer Interaction* 2, CSCW (Nov. 2018), 1–29. <https://doi.org/10.1145/3274374>
- [52] Anthony Liu, Santiago Guerra, Isaac Fung, Gabriel Matute, Ece Kamar, and Walter Lasecki. 2020. Towards Hybrid Human-AI Workflows for Unknown Unknown Detection. In *Proceedings of The Web Conference 2020*. ACM, Taipei Taiwan, 2432–2442. <https://doi.org/10.1145/3366423.3380306>
- [53] Christina Low, Emma McCamey, Cole Gleason, Patrick Carrington, Jeffrey P. Bigham, and Amy Pavel. 2019. Twitter A11y: A Browser Extension to Describe Images. In *The 21st International ACM SIGACCESS Conference on Computers and Accessibility*. ACM, Pittsburgh PA USA, 551–553. <https://doi.org/10.1145/3308561.3354629>
- [54] Michael A. Madaio, Luke Stark, Jennifer Wortman Vaughan, and Hanna Wallach. 2020. Co-Designing Checklists to Understand Organizational Challenges and Opportunities around Fairness in AI. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. ACM, Honolulu HI USA, 1–14. <https://doi.org/10.1145/3313831.3376445>
- [55] Sally Maitlis and Marlys Christianson. 2014. Sensemaking in Organizations: Taking Stock and Moving Forward. *Academy of Management Annals* 8, 1 (Jan. 2014), 57–125. <https://doi.org/10.5465/19416520.2014.873177>
- [56] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model Cards for Model Reporting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*. ACM, Atlanta GA USA, 220–229. <https://doi.org/10.1145/3287560.3287596>
- [57] Michael Muller, Ingrid Lange, Dakuo Wang, David Piorkowski, Jason Tsay, Q. Vera Liao, Casey Dugan, and Thomas Erickson. 2019. How Data Science Workers Work with Data: Discovery, Capture, Curation, Design, Creation. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM, Glasgow Scotland Uk, 1–15. <https://doi.org/10.1145/3290605.3300356>
- [58] Elizamary de Souza Nascimento, Iftekhar Ahmed, Edson Oliveira, Marcio Piedade Palheta, Igor Steinmacher, and Tayana Conte. 2019. Understanding Development Process of Machine Learning Systems: Challenges and Solutions. In *2019 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM)*. IEEE, Porto de Galinhas, Recife, Brazil, 1–6. <https://doi.org/10.1109/ESEM.2019.8870157>
- [59] NTSB. 2018. Preliminary Report HWY18MH010. (2018), 4.
- [60] Besmira Nushi, Ece Kamar, and Eric Horvitz. 2018. Towards Accountable AI: Hybrid Human-Machine Analyses for Characterizing System Failure. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, Vol. 6. 10.
- [61] Luke Oakden-Rayner, Jared Dunnmon, Gustavo Carneiro, and Christopher Re. 2020. Hidden stratification causes clinically meaningful failures in machine learning for medical imaging. In *Proceedings of the ACM Conference on Health, Inference, and Learning*. ACM, Toronto Ontario Canada, 151–159. <https://doi.org/10.1145/3368555.3384468>

- [62] David Piorkowski, Soya Park, April Yi Wang, Dakuo Wang, Michael Muller, and Felix Portnoy. 2021. How AI Developers Overcome Communication Challenges in a Multidisciplinary Team: A Case Study. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW1 (April 2021), 1–25. <https://doi.org/10.1145/3449205>
- [63] Peter Pirolli and Stuart Card. 2005. The Sensemaking Process and Leverage Points for Analyst Technology as Identified Through Cognitive Task Analysis. *Proceedings of International Conference on Intelligence Analysis* 2005, January (2005), 2–4. <https://doi.org/10.1007/s13398-014-0173-7.2>
- [64] Mahima Pushkarna, James Wexler, and Jimbo Wilson. 2017. Facets: An Open Source Visualization Tool for Machine Learning Training Data. <https://pair-code.github.io/facets/>
- [65] Iyad Rahwan, Manuel Cebrian, Nick Obradovich, Josh Bongard, Jean-François Bonnefon, Cynthia Breazeal, Jacob W. Crandall, Nicholas A. Christakis, Iain D. Couzin, Matthew O. Jackson, Nicholas R. Jennings, Ece Kamar, Isabel M. Kloumann, Hugo Larochelle, David Lazer, Richard McElreath, Alan Mislove, David C. Parkes, Alex ‘Sandy’ Pentland, Margaret E. Roberts, Azim Shariff, Joshua B. Tenenbaum, and Michael Wellman. 2019. Machine Behaviour. *Nature* 568, 7753 (April 2019), 477–486. <https://doi.org/10.1038/s41586-019-1138-y>
- [66] Alexander Ratner, Stephen H. Bach, Henry Ehrenberg, Jason Fries, Sen Wu, and Christopher Ré. 2017. Snorkel: rapid training data creation with weak supervision. *Proceedings of the VLDB Endowment* 11, 3 (Nov. 2017), 269–282. <https://doi.org/10.14778/3157794.3157797>
- [67] Donghao Ren, Saleema Amershi, Bongshin Lee, Jina Suh, and Jason D. Williams. 2017. Squares: Supporting Interactive Performance Analysis for Multiclass Classifiers. *IEEE Transactions on Visualization and Computer Graphics* 23, 1 (Jan. 2017), 61–70. <https://doi.org/10.1109/TVCG.2016.2598828>
- [68] Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond Accuracy: Behavioral Testing of {NLP} Models with {C}heck{L}ist. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 4902–4912. <https://doi.org/10.18653/v1/2020.acl-main.442>
- [69] George Robertson, Mary Czerwinski, Kevin Larson, Daniel C. Robbins, David Thiel, and Maarten van Dantzich. 1998. Data Mountain: Using Spatial Memory for Document Management. In *Proceedings of the 11th annual ACM symposium on User interface software and technology - UIST '98*. ACM Press, San Francisco, California, United States, 153–162. <https://doi.org/10.1145/288392.288596>
- [70] Yvonne. Rogers. 2012. *HCI Theory*.
- [71] Daniel M. Russell, Mark J. Stefik, Peter Pirolli, and Stuart K. Card. 1993. The cost structure of sensemaking. In *Proceedings of the SIGCHI conference on Human factors in computing systems - CHI '93*. ACM Press, Amsterdam, The Netherlands, 269–276. <https://doi.org/10.1145/169059.169209>
- [72] Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh, and Lora M Aroyo. 2021. “Everyone wants to do the model work, not the data work”: Data Cascades in High-Stakes AI. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. ACM, Yokohama Japan, 1–15. <https://doi.org/10.1145/3411764.3445518>
- [73] Alex Serban, Koen van der Blom, Holger Hoos, and Joost Visser. 2020. Adoption and Effects of Software Engineering Best Practices in Machine Learning. In *Proceedings of the 14th ACM / IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM)*. ACM, Bari Italy, 1–12. <https://doi.org/10.1145/3382494.3410681>
- [74] Helen Shen. 2014. Interactive notebooks: Sharing the code. *Nature* 515, 7525 (Nov. 2014), 151–152. <https://doi.org/10.1038/515151a>
- [75] Hong Shen, Haojian Jin, Ángel Alexander Cabrera, Adam Perer, Haiyi Zhu, and Jason I. Hong. 2020. Designing Alternative Representations of Confusion Matrices to Support Non-Expert Public Understanding of Algorithm Performance. *Proceedings of the ACM on Human-Computer Interaction* 4, CSCW2 (Oct. 2020), 1–22. <https://doi.org/10.1145/3415224>
- [76] Connor Shorten and Taghi M. Khoshgoftaar. 2019. A Survey on Image Data Augmentation for Deep Learning. *Journal of Big Data* 6, 1 (Dec. 2019), 60. <https://doi.org/10.1186/s40537-019-0197-0>
- [77] Sahil Singla, Besmira Nushi, Shital Shah, Ece Kamar, and Eric Horvitz. 2021. Understanding Failures of Deep Networks via Robust Feature Extraction. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021*. Computer Vision Foundation / IEEE. https://openaccess.thecvf.com/content/CVPR2021/papers/Singla_Understanding_Failures_of_Deep_Networks_via_Robust_Feature_Extraction_CVPR_2021_paper.pdf
- [78] Anselm L. Strauss and Juliet M. Corbin (Eds.). 1997. *Grounded theory in practice*. Sage Publications, Thousand Oaks.
- [79] Yuchi Tian, Kexin Pei, Suman Jana, and Baishakhi Ray. 2018. DeepTest: Automated Testing of Deep-Neural-Network-Driven Autonomous Cars. In *Proceedings of the 40th International Conference on Software Engineering*. ACM, Gothenburg Sweden, 303–314. <https://doi.org/10.1145/3180155.3180220>
- [80] Stefan Timmermans and Iddo Tavory. 2012. Theory Construction in Qualitative Research: From Grounded Theory to Abductive Analysis. *Sociological Theory* 30, 3 (Sept. 2012), 167–186. <https://doi.org/10.1177/0735275112457914>
- [81] Deniz Ucbasaran, Dean A. Shepherd, Andy Lockett, and S. John Lyon. 2013. Life After Business Failure: The Process and Consequences of Business Failure for Entrepreneurs. *Journal of Management* 39, 1 (Jan. 2013), 163–202. <https://doi.org/10.1177/0149206312457914>

[//doi.org/10.1177/0149206312457823](https://doi.org/10.1177/0149206312457823)

- [82] Zhiyuan Wan, Xin Xia, David Lo, and Gail C. Murphy. 2020. How does Machine Learning Change Software Development Practices? *IEEE Transactions on Software Engineering* (2020), 1–1. <https://doi.org/10.1109/TSE.2019.2937083>
- [83] Haohan Wang, Xindi Wu, Zeyi Huang, and Eric P. Xing. 2020. High-Frequency Component Helps Explain the Generalization of Convolutional Neural Networks. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Seattle, WA, USA, 8681–8691. <https://doi.org/10.1109/CVPR42600.2020.00871>
- [84] Taowei David Wang, Krist Wongsuphasawat, Catherine Plaisant, and Ben Shneiderman. 2011. Extracting Insights from Electronic Health Records: Case Studies, a Visual Analytics Process Model, and Design Recommendations. *Journal of Medical Systems* 35, 5 (Oct. 2011), 1135–1152. <https://doi.org/10.1007/s10916-011-9718-x>
- [85] Karl E. Weick. 1995. *Sensemaking in Organizations*. Sage Publications, Thousand Oaks.
- [86] Karl E. Weick, Kathleen M. Sutcliffe, and David Obstfeld. 2005. Organizing and the Process of Sensemaking. *Organization Science* 16, 4 (Aug. 2005), 409–421. <https://doi.org/10.1287/orsc.1050.0133>
- [87] James Wexler, Mahima Pushkarna, Tolga Bolukbasi, Martin Wattenberg, Fernanda Viegas, and Jimbo Wilson. 2019. The What-If Tool: Interactive Probing of Machine Learning Models. *IEEE Transactions on Visualization and Computer Graphics* (2019), 1–1. <https://doi.org/10.1109/TVCG.2019.2934619>
- [88] Gail Whiteman and William H. Cooper. 2011. Ecological Sensemaking. *Academy of Management Journal* 54, 5 (Oct. 2011), 889–911. <https://doi.org/10.5465/amj.2008.0843>
- [89] Tongshuang Wu, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel Weld. 2019. {E}rrudite: Scalable, Reproducible, and Testable Error Analysis. *Proceedings of the 57th Conference of the Association for Computational Linguistics* (2019), 747–763. <https://www.aclweb.org/anthology/P19-1073>
- [90] Qian Yang, Jina Suh, Nan-Chen Chen, and Gonzalo Ramos. 2018. Grounding Interactive Machine Learning Tool Design in How Non-Experts Actually Build Models. In *Proceedings of the 2018 Designing Interactive Systems Conference*. ACM, Hong Kong China, 573–584. <https://doi.org/10.1145/3196709.3196729>
- [91] Jie M. Zhang, Mark Harman, Lei Ma, and Yang Liu. 2020. Machine Learning Testing: Survey, Landscapes and Horizons. *IEEE Transactions on Software Engineering* (2020), 1–1. <https://doi.org/10.1109/TSE.2019.2962027>
- [92] Mengshi Zhang, Yuqun Zhang, Lingming Zhang, Cong Liu, and Sarfraz Khurshid. 2018. DeepRoad: GAN-based metamorphic testing and input validation framework for autonomous driving systems. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*. ACM, Montpellier France, 132–142. <https://doi.org/10.1145/3238147.3238187>
- [93] Xufan Zhang, Yilin Yang, Yang Feng, and Zhenyu Chen. 2019. Software Engineering Practice in the Development of Deep Learning Applications. *arXiv:1910.03156 [cs]* (Oct. 2019). <http://arxiv.org/abs/1910.03156>