# MTV: Visual Analytics for Detecting, Investigating, and Annotating Anomalies in Multivariate Time Series

DONGYU LIU, Massachusetts Institute of Technology

SARAH ALNEGHEIMISH, Massachusetts Institute of Technology

ALEXANDRA ZYTEK, Massachusetts Institute of Technology

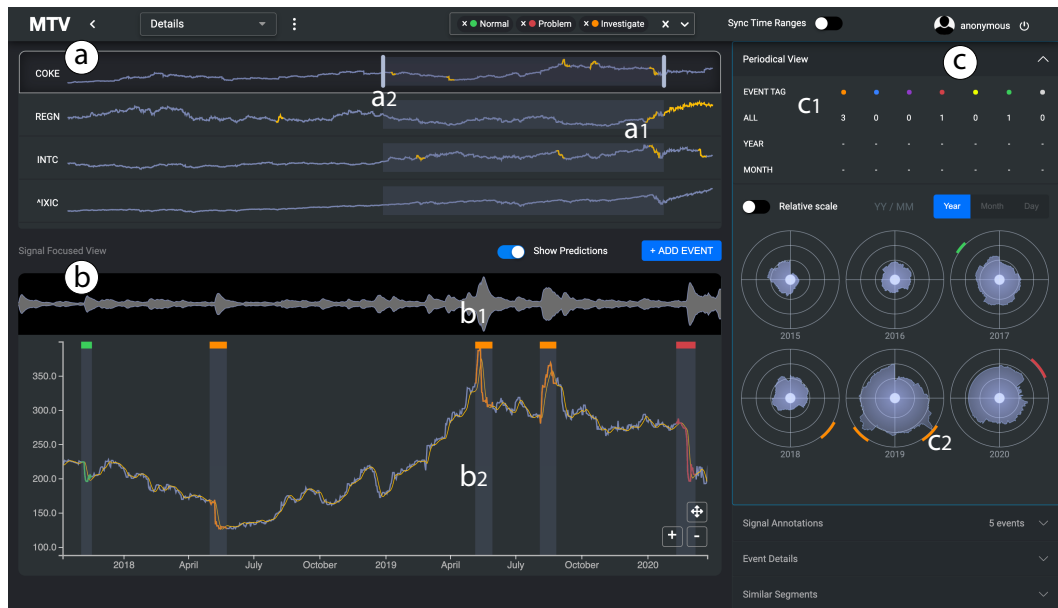KALYAN VEERAMACHANENI, Massachusetts Institute of Technology

Fig. 1. The MTV interface, displaying an analysis of stock data. The Signal Overview (a) displays multiple signals (in this case, stocks) that share the same horizontal timeline and highlights anomalous events with a warning color (a1). Users pick a signal of interest (i.e., COKE) and brush a segment (a2) to observe its details and interact with the events in the Signal Focused View (b). Events can be color-tagged (e.g., green - normal, orange - investigate, red - problem) and filtered in the top header. The predicted errors (b1) can be toggled to explain why a certain event (b2) was identified by the machine learning algorithm. The Side Panel (c) includes four collapsible views — the Periodical View (c1, c2), Signal Annotations View, Event Details View, and Similar Segments View — which allow users to investigate and annotate events efficiently and collaboratively.

Detecting anomalies in time-varying multivariate data is crucial in various industries for the predictive maintenance of equipment. Numerous machine learning (ML) algorithms have been proposed to support

automated anomaly identification. However, a significant amount of human knowledge is still required to interpret, analyze, and calibrate the results of automated analysis. This paper investigates current practices used to detect and investigate anomalies in time series data in industrial contexts and identifies corresponding needs. Through iterative design and working with nine experts from two industry domains (aerospace and energy), we characterize six design elements required for a successful visualization system that supports effective detection, investigation, and annotation of time series anomalies. We summarize an ideal human-AI collaboration workflow that streamlines the process and supports efficient and collaborative analysis. We introduce MTV (Multivariate Time Series Visualization), a visual analytics system to support such workflow. The system incorporates a set of novel visualization and interaction designs to support multi-faceted time series exploration, efficient in-situ anomaly annotation, and insight communication. Two user studies, one with 6 spacecraft experts (with routine anomaly analysis tasks) and one with 25 general end-users (without such tasks), are conducted to demonstrate the effectiveness and usefulness of MTV.

CCS Concepts: • **Human-centered computing** → **Visual analytics**.

Additional Key Words and Phrases: anomaly analysis, time series, visual analytics, collaborative analysis, annotation, human-AI collaboration

## 1 INTRODUCTION

The rapid proliferation of sensors and connected devices has led to a massive, ever-increasing accumulation of temporal observation data (a.k.a. multivariate time series). As more and more industries rely on such data, detecting and analyzing anomalies in time series becomes increasingly important for critical use cases, ranging from cyber-intrusion and fault detection to preventative maintenance and fraud prevention [1].

A time series anomaly is a time point or period during which a system exhibits unusual behavior. To efficiently detect such unusual behaviors in what are often enormous datasets, and ease the burden on the humans often tasked with such detection, a variety of automated anomaly detection methods using statistical or machine learning techniques have been proposed [16, 23, 25, 32]. A typical approach, employed by applied machine learning engineers, is to use an ML technique (possibly a deep learning approach), flag anomalous periods or time points, and send them to experts/operators, often using simple csv files.

In practice, however, these methods only solve part of the problem — detecting unusual behavior. Many real-world systems are highly dynamic, and can only be fully understood by including context and domain expertise outside the scope of the time series analyzed by these methods. Simply put, the presence of unusual behavior alone does not necessarily mean there is a *problem* that needs attention. For example, if a satellite is passing an eclipse, its signals (one of the datasets we consider) might exhibit some unusual behaviors that would not necessarily indicate problems. Additionally, system dynamics, like maneuvers, mode switching and others, may manifest in signals as unusual but are not necessarily troublesome. In our collaborations with domain experts and operators, we found that their reluctance to rely on the ML-based tactics above stemmed from their inability to integrate their own knowledge, context and analysis with the ML approach, to share their analyses with one another, to understand the outputs of ML approaches, and ultimately to reduce the number of false positives over time. Thus, we set out to study and design a flexible human-AI collaboration workflow to overcome this reluctance.

Data visualization supports this workflow by integrating humans' valuable expertise and creativity into anomaly investigation via various visual representations of information and human-machine

interactions [68]. Recent years have seen increasing efforts to combine visualization and automated anomaly detection for time series data analysis [72, 74]. However, effectively visualizing multivariate time series is difficult, as it requires determining how to get across complex yet central aspects of the data — including issues of scalability (what is the right level of detail to show), multi-dimensionality (how many time series should be displayed at once), and interpretability (how best to represent the contextual information that helps determine legitimate anomalies).

Other questions come up as well: (1) **End-user:** A large number of end-users tasked with time series anomaly analysis may have minimal or no technical background. To help them, we must understand what they are hoping to _observe_ or _control_ about the data and the analysis process, and avoid using overly technical jargon or complicated visualizations. An intuitive visual design is required to alleviate the learning burden for users. (2) **Annotation:** From the user's point of view, it is the interpretation of results that matters. However, existing visual systems [14, 73, 74] for time series anomaly analysis lack the ability to document and share such interpretations, without which it can be difficult to translate results into actionable decisions. We stress that the ability to annotate anomalies, as well as to organize and present these annotations, is a highly important part of any anomaly detection system. (3) **Workflow:** Bringing humans into the analysis loop to make sense of ML results is not a simple or one-step task. Still more questions arise during this process, including: What ML outputs does it make sense to start with? How should teams collaborate to annotate anomalies? Can the ML system use these annotations to reduce false alarms in the future? Little existing work has investigated these questions deeply.

Awareness of these gaps comes through our collaboration with domain experts in the aerospace and energy industries, who spend a considerable amount of time monitoring telemetry data from large industrial devices in order to identify and analyze potentially hazardous events. This work represents an early step in understanding the support-annotating and sensemaking needs of these experts — who, like many in industrial environments, work with a large number of large-scale time series, each of which may contain more than tens of thousands of data points. In this context, we lay out research questions and summarize our key contributions as follows. The results of RQ1, RQ2, and RQ3 are all derived from an iterative user-centered design process which involves 6 experts from aerospace and 3 experts from energy industry.

- **RQ1.** _What design elements are required for a system to support time series anomaly analysis in the above-mentioned context?_
  We summarize the current challenges commonly faced by both the aerospace and energy experts (Section 3.2). Then we identify the design requirements (Section 3.3) necessary for a time series analysis system that incorporates both the power of human intelligence and the computation capability of machines to solve existing challenges.
- **RQ2.** _What is an ideal human-AI collaboration workflow for efficient and effective detection, investigation, and annotation of anomalies in industrial-scale time series data?_
  We summarize a streamlined human-AI collaboration workflow (Section 3.4) that allows for easy, flexible, and efficient anomaly analysis of multiple large-scale time series data. We report the usability challenges involved in such a workflow, and discuss the lessons learned that also apply to other human-in-the-loop data analytics scenarios, in Section 6.
- **RQ3.** _What is an effective system solution that enables such human-AI collaboration workflow while incorporating all the design elements learned from RQ1?_
  We develop MTV, a visual analytics system that follows the aforementioned design requirements and demonstrates our proposed workflow. The system includes an end-to-end ML pipeline that detects anomalies without labeled data (Section 4.1) and learns from human feedback (Section 4.2). The interface (Section 4.4) incorporates several novel and intuitive

visualization and interaction designs for multi-faceted and variously granular time series data exploration, as well as anomaly investigation and in-situ annotation and communication. Specifically, a novel shape-matching algorithm (Section 4.3) is introduced to enhance annotation efficiency. We quantitatively evaluate the shape-matching algorithm, as well as the system's ability to learn from human feedback, and report the results in Section 5.3.

- **RQ4.** *How do domain experts, who are well-versed in routine anomaly analysis tasks, perceive the usefulness and usability of such a system? How about general end-users with less experience?* We evaluate the usefulness of MTV through several case studies, performed by 6 domain experts from the aerospace industry and using spacecraft telemetry data (Section 5.1). We also run a usability study using stock data in which 25 general end-users reveal the potential broad benefits of such a system (Section 5.2).

## 2 RELATED WORK

In this section, we summarize the techniques that are most relevant to our work. We have divided these into four types: automated anomaly detection, visualization for time series anomaly detection, annotating with interactive visualizations, and sensemaking and collaborative visualization.

### 2.1 Automated Anomaly Detection

Over the decades, the rich variety of anomaly types, data types and application scenarios has led to the development of numerous anomaly detection approaches. Several surveys have summarized these techniques [16, 23, 25, 32].

The general goal of anomaly detection is to find unexpected patterns in data. The simplest approaches are "out-of-limits" methods, which are applied to raw values and flag locations where predefined thresholds are surpassed. However, such methods are not suited for detecting *contextual anomalies* that do not fall within low-density values in global but are anomalous with respect to local values. To overcome this, advanced approaches have since been developed based on statistics [54, 54, 76], clustering [6, 13], and machine learning [1, 5] or deep learning [21, 36, 77].

In our work, we build an end-to-end ML pipeline which integrates three powerful algorithms including Arima, LSTM and TadGAN to handle datasets of different characteristics (Section 4.1). Our main goal is to engage humans with the analysis loop, enabling them to make sense of the ML results (especially the identified contextual anomalies), and to support collaborative sensemaking between human experts.

### 2.2 Visualization for Time Series Anomaly Detection

There are numerous techniques for visualizing and structuring time series [2]. The key difference lies in how the timeline is encoded. Time is linear, but contains an inherent hierarchical structure of granularities, such as hours, days, weeks, and months. Depending on the analytical tasks at hand, different design techniques may be employed.

A standard method for visualizing a time series is mapping time to the horizontal x-axis and time-dependent variables to the vertical y-axis [2]. Line and area charts are the most common ways to represent time series. When our focus is on observing cyclic or periodical patterns, a spiral-shaped time axis [69] is a useful time-encoding scheme. If we put more emphasis on individual dates, a calendar layout is more suitable, and can depict daily, monthly, or yearly value changes [64].

To encode multivariate time series, techniques including superposed line graphs, braided graphs, small multiples, and horizon graphs can be used [42]. Recently, glyph-based designs have also been explored due to their expressiveness and effective use of screen space [15, 19]. Another typical approach involves using multiple coordinated views, with each view displaying particular aspects of time in order to support a coordinated analysis [2]. Each of the techniques introduced above has

its own merits and is well-suited for certain analytical tasks. However, it is difficult to use any one of them for anomaly analysis directly.

Several comprehensive visual analytics systems have been developed for time series anomaly investigation. These systems have been applied to various scenarios, such as air quality monitoring [55], traffic volume monitoring [14], electronic healthcare records analysis [73], and cloud computing system performance analysis [74]. None of them support anomaly annotation and collaborative analysis.

In this work, we propose a set of hybrid visualization and interaction designs tailored for time series anomaly analysis with a particular focus on collaborative annotations. All these designs together compose our MTV system, which supports a streamlined and effective human-AI collaboration workflow (Section 3.4).

## 2.3   Annotating with Interactive Visualizations

Annotating (labeling and commenting on) data is frequently supported in computer-supported cooperative work and human-computer interaction (CSCW/HCI) [7, 8, 10, 51]. Annotation tasks are widely observed in many application scenarios, as humans leverage their high-level intelligence and domain expertise to add meaningful context to data. Through visual interfaces, humans can annotate interesting entities in textual documents [31, 71], images [9], and videos [33]. Unlike text and image data that are often easily understood by humans, time series data are much harder to annotate, and studies regarding time series data annotation are scarce [3, 57]. To the best of our knowledge, no existing work focuses on large-scale multivariate time series data annotation through the lens of human-machine collaboration and collaborative analysis among human experts.

Combining machine learning methods and visualization techniques can significantly improve annotation efficacy [31, 33, 59]. Unsupervised techniques are often employed to ease the burden of annotation tasks by recommending important instances to annotate. One typical approach is to use dimensionality reduction techniques to plot instances in the 2D plane, along with instance selection interactions, to allow for easy annotation [8]. Our goal is to investigate time series anomalies and create meaningful annotations for them, which requires having contextual temporal information for anomalous segments. Thus, we propose the use of multiple coordinated views, with each view showing different contextual temporal information, to support in-situ annotation and communication. We further present several novel techniques, including a Multi-Aggregation Viewer and a Similar Segments View powered by shape-matching algorithms, to facilitate anomaly annotation and boost efficiency.

## 2.4   Sensemaking and Collaborative Visualization

Collaborative visualization is defined as "*the shared use of computer-supported, (interactive) visual representations of data by more than one person with the common goal of contribution to joint information processing activities*" by Isenberg et al. [38]. Prior research in CSCW/HCI has demonstrated the importance of collaborative visualization for information-seeking and decision-making [4, 34, 35, 39, 46, 52, 53, 60], and for sensemaking of a variety types of application data such as tweet data [63], mobile data [47], sound data [18] and crime data [24]

Sensemaking in collaborative visual analytics is challenging. Experts must have iterative discussions in order to form and verify hypotheses, draw conclusions, and publish findings [49]. Members of a team must also maintain *awareness* of each other's work in order to progress [20]. (In our scenario, the "work" refers to annotations on anomalous events.) Many existing tools enable team members to record questions and insights — in the form of text, audio, or visual diagrams — to facilitate organizing and sharing their results [12, 26, 29, 40, 49, 65, 70].

Compared with synchronous collaboration (through shared workplaces or real-time networked displays) [37, 39, 40, 49], remote asynchronous collaboration [28, 30, 75] remains relatively unexplored, particularly in the context of time series data. Compared with text or images, time series are intrinsically difficult for humans to interpret without sufficient contextual information. This makes it even harder for annotators to create accurate externalizations (external representations of a person's internal thoughts) in order to communicate. How to encode and display these externalizations, and how to maintain awareness between team members in a remote asynchronous setting, are still open research questions. In this work, we summarize a streamlined human-AI collaboration workflow that supports efficient collaborative anomaly annotations. More importantly, we propose a set of novel visualizations and interactions to support in-situ annotation and communication, as well as to enhance team members' awareness of each other's findings.

## 3 METHOD

### 3.1 Iterative Design

We followed an iterative user-centered design process to develop MTV. We collaborated with industrial domain experts, gathering design requirements and collecting feedback from them. We began with proof-of-concept mockup designs made using Figma (an online collaborative interface design tool). We then moved on to an interactive system prototype, a high-fidelity prototype, and eventually to a deployable system.

*3.1.1 **Participants:*** The entire design process involved 9 participants in two groups.

The primary group consisted of 6 domain experts from a world-leading communication satellite company: one spacecraft program manager (P1) and five senior satellite engineers (P2-P6). The team analyzes spacecraft telemetry data for signs of hazards that may result in system failure. Each expert has between 5 and 17 years of spacecraft telemetry data analysis experience, and between 0 and 3 years of machine learning experience. This group participated in our regular design meetings, which occurred once or twice per month for three years, and provided feedback on prototypes.

The external group included 3 senior engineers (E1-E3) from an energy company with expertise in analyzing time series from wind turbines and energy pipelines. They took surveys about the system design and usability, and tested MTV with their own data. This group provided assurance that the design requirements and features we developed with the primary group are generalizable in certain contexts (defined in Fig. 2), as well as applicable to similar fields.

*3.1.2 **Design stages:*** The design process occurred in three stages.

We began the first stage by conducting interviews with the experts from the primary group (P1-P6) to understand in detail the methods they use to analyze spacecraft telemetry data for anomalies as part of their routine workload, along with the associated challenges. We detail this process in Section 3.2. Knowing all these challenges, we had additional conversions with E1-E3 and confirmed that the same challenges also exist in their domain.

In this first stage, we gathered initial requirements and designed mockups using Figma by choosing different combinations of visualization techniques in accordance with design principles. We refined the requirements iteratively, and finally derived a set of requirements to guide our development of the early version of MTV.

In the second stage, we focused on implementing MTV and iteratively improving it by engaging domain experts from the primary group. We then ran our high-fidelity prototype on public data. This process involved many informal discussions and interviews. From time to time, we also presented the prototype to E1-E3 and collected feedback from them. The output of this stage was a high-fidelity software prototype.

Lastly, we deployed our prototype using real spacecraft telemetry data in a production environment, and conducted four case studies with P1-P6 through online meetings with screen sharing. Based on the feedback collected in each case study, we continued to refine our system. The experiment details and results are reported in Section 5.1.

## 3.2 How Satellite Experts Detect and Investigate Anomalies in Spacecraft Telemetry?

One major objective of the satellite company's operative team is to detect unexpected behaviors (i.e., anomalies) in tens of thousands of signals. Through our three-year collaboration with the satellite company, we formed an understanding of the scope of this task and what it entails. The team works with multiple spacecrafts. Each spacecraft telemetry database contains around 37,000 signals spanning 9 different subsystems from one spacecraft. Each signal is a univariate time series collected at the microsecond level, and has been tracked for over 10 years.

The team's conventional approach to anomaly detection is based around setting and adjusting thresholds in order to flag anomalous intervals. The team then manually reviews suspicious intervals, often using simple `csv` files, and examines individual signals in a third-party platform such as `MatLAB`. For each anomaly detected, the team digs into the corresponding signal to track the root cause of the alarm. During this investigation process, some relevant signals (usually from the same subsystem) will be examined. An average of 20 alarms is reported every day, most of which are false alarms and can be resolved within a few hours. For some that are identified as true alarms but non-urgent, the experts gather further information over some time window to help explain the root cause and the way forward.

Over the course of this process, the team faced many challenges:

**C1** The team's current method is highly expensive, and forces them to restrict their focus to a subset of a few hundred signals chosen based on domain knowledge.

**C2** The team has wanted to use ML models to identify contextual anomalies — anomalies that do not exceed a normal range, but are unusual compared to local values. But they have limited machine learning experience and have not been able to find a way to do this.

**C3** The team found that ML models often flag unusual patterns, even when these patterns do not necessarily indicate a problem. For example, an eclipse might cause patterns that are then flagged even though they are not troublesome. The team is eager for their models not to mark these patterns, but struggles to integrate this feedback.

**C4** To track the root cause of an anomalous event, the team often needs to observe several relevant signals. As such, they have to frequently shift between different tools such as `Excel` and `MatLab`. This process becomes particularly inefficient when the number of signals that need to be analyzed together increases.

**C5** Frequently, the team needs informal discussion between team members. However, it is arduous for the team to document what they found and communicate their insights effectively, which is the same challenge they face during the in-person meeting.

## 3.3 Design Requirements

Motivated by the aforementioned challenges, we came up with the following design requirements, which guided the development of MTV.

**R1 Facilitate efficient and intelligent anomaly identification (C1, C3).**
The scale of time series data is often overwhelming in terms of either the number of signals or the number of data points per signal. The system should support fast and robust automated anomaly detection. In addition, the system is expected to learn from human annotations (e.g.,

tags and comments) of existing anomalous events in order to avoid generating too many false alarms or missing true anomalies.

**R2 Provide the flexibility to interact with ML results (C2).**
ML can exhaust every meaningful hyperparameter combination suggested by experts. Different settings may lead to different results that make sense in certain analytical scenarios. For example, the satellite experts may want to analyze ML results over different time aggregation levels, such as "6 minutes" and "1 hour" (Section 4.1). The system should save possibly meaningful ML results and allow experts to interact with these results in real time to choose what they want to investigate.

**R3 Offer a visual interface to display and interact with ML-flagged events (C3, C4).**
The visual interface should display when and where the anomalies occurred and point experts to the most important time segments, prioritized by severity. Create, read (aka retrieve), update, and delete (CRUD) functions should be supported when experts interact with an event. For example, when an expert finds a segment important but it is not flagged by the machine, s/he should be allowed to easily create an event on this segment.

**R4 Allow efficient in-situ annotation and communication (C5).**
The system should allow experts to document their thoughts, as well as to assign tags to an event. We refer to this behavior as event annotation. In-situ annotation [7] — in our scenario, annotating directly on the time segment of the event under investigation — should be supported. The system should also allow in-situ communication of insights (i.e., tags and textual comments) and enhance team members' awareness of each other's findings. This helps experts to make well-informed annotations, and to verify each other's decisions with their complementary knowledge. Techniques should be used to assist with efficient annotation, as the large-scale nature of the data may lead to too many anomalies being identified.

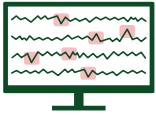**R5 Enable multi-faceted and multi-granular visual explorations of time series (C4).**
To interpret and annotate anomalies, the system should provide experts with contextual information about the anomalies, allowing them to confirm or refute the system's conclusions. Therefore, the ability to show time series patterns from different perspectives (trend, periodicity, etc.) and in different granularities (second, minute, hour, etc) is required.

**R6 Afford a streamlined workflow for efficient collaboration (C1-C5).**
Experts are expected to follow a linear workflow to perform anomaly analysis. First, the workflow should meet all the previous requirements (**R1-R5**). Second, experts are allowed to enter the workflow at any stage in order to do their job without starting from scratch. For example, expert A may specify what ML results to investigate, while experts B and C perform detailed analysis and make annotations collaboratively and expert D engages in discussion around an anomaly; expert A may then come back to check the final annotation results. In short, the pathway to any individual task should be directly accessible.

## 3.4 Workflow

Through iterations, we identify an ideal human-AI collaboration workflow (**R6**) that streamlines the process of time series anomaly analysis. Experts can enter the workflow at any stage and finish their tasks in an efficient and collaborative manner. Table 1 summarizes all the relevant steps in detail. The workflow starts with ML extracting anomalies from massive time series, continues with users giving high-level input to machines about what to observe, proceeds with an overview-first and details-on-demand exploration (a.k.a. Shneiderman's mantra [61]) and a fast in-situ annotation strategy, and ends with enhancing the ML pipeline through annotations, forming a closed loop.
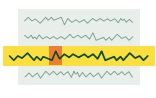
### T1. Extract anomalies from massive time series (R1). [ML]

The first step is to use ML to detect anomalous events within our dataset, which consists of multiple (more than 100) large-scale time series (often more than 10K data points/time steps) without labeled data. An end-to-end unsupervised ML pipeline is built for this purpose. The pipeline exposes analytical-task-relevant hyperparameters to experts, such as the time interval for aggregating raw signals and strategies to impute missing values. The ML results of every meaningful hyperparameter setting are saved.
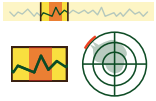
### T2. Select signals and ML results to investigate (R2). [Human]

Depending on need, an expert selects a subset of signals (a few to dozens) and ML pipeline outputs to investigate. For example, s/he may want to look at anomalies identified in all temperature signals (~20 of them) aggregated at an one-hour level.
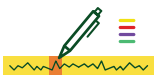
### T3. Scan through these signals and pick one of interest (R3, R5). [Human]

After selecting the signals, the expert may want to observe an overview of them in which multiple time series are plotted and anomalous events are highlighted. This allows s/he to quickly scan through the dynamics of multiple signals and then decide which to pick for further investigation.

### T4. Explore one signal and its anomalies in depth (R4, R5). [Human]

Once the expert finds the signal and anomalies of interest, s/he may explore contextual information about the anomalies. This information includes but is not limited to: (1) the surrounding time series values of the anomaly at different aggregation levels; (2) the periodical (daily/monthly/yearly) patterns of the time series; (3) any annotations from other team members; and (4) an anomaly's total interaction/annotation history.

### T5. Annotate the anomalies from this time series (R3, R4). [ML+Human]

After gaining a sufficient understanding of one anomalous event, the expert may interact with (CRUD) or annotate (tag and comment) the event. If this event has already been annotated by other team members, the expert can still add more comments and change the tag if necessary. Given the large-scale nature of the data, it is challenging for experts to annotate all events. Automatic techniques such as shape-matching are expected to boost the efficiency of annotation.

### T6. Enhance ML pipeline using annotations (R1). [ML+Human]

All annotations and CRUD interactions are stored in the database. These are leveraged to improve the performance of the ML pipeline, forming a closed loop. In other words, we incorporate experts' knowledge into the system so that it can learn from event patterns, mitigating false alarms in favor of surfacing true anomalies in the future.

Table 1. The ideal human-AI collaboration workflow that streamlines tasks. It starts with ML extracting anomalies from massive time series, continues with users giving high-level input to machines about what to observe, proceeds with an overview-first and details-on-demand exploration and a fast in-situ annotation strategy, and ends with enhancing the ML pipeline through annotations, forming a closed loop.

| Decision Risk | | Decision Time | | Domain Expertise | | Technical Expertise | | Role of ML | |
|---|---|---|---|---|---|---|---|---|---|
| Low | ✓ | A few seconds | | Low | | Low | ✓ | Clarify understanding | |
| Medium | ✓ | Less than a minute | | Medium (intuition) | ✓ | Medium | ✓ | Improve trust | |
| High | | A few minutes | | High | ✓ | High | ✓ | Handle disagreements | ✓ |
| | | A few hours | ✓ | | | | | Accountability | |
| | | A few days | ✓ | | | | | Expedite decisions | ✓ |
| | | | | | | | | Learn about domain | ✓ |
| | | | | | | | | Curiosity | |

✓ The context factors that apply to MTV

Fig. 2. The five domain context factors we used to determine the scenarios in which our proposed workflow and system are a good fit for time series anomaly analysis. *Decision* refers to the process through which a team judges an ML-identified anomaly and provides a corresponding explanation. MTV is designed to support asynchronous collaborative analysis, and is thereby suited for decisions of that are of low or medium risk and take at least a few hours. The system is friendly to non-technical users but requires a certain amount of domain expertise so that time series can be interpreted.

It is worth noting that this workflow is not meant to fit every single analytical scenario. To clarify its scope, we consider five domain context factors (shown in Fig. 2) inspired by the work of Zytek et. al. [78]. To instantiate this workflow, we developed MTV, an interactive visual analytics system for time series anomaly analysis at an industrial scale. The system supports the foregoing human-AI collaboration workflow and integrates all the features to meet the design requirements.

## 4 MTV

In this section, we first describe our approach of identifying anomalies through an end-to-end machine learning pipeline, followed by the introduction of improving ML with annotations. Next, we detail how we use shape-matching techniques to enhance annotation efficiency. Finally, we describe the visual design of our system.

### 4.1 Identify Anomalies with the End-to-end ML Pipeline

| Signal ID | Event ID | $t_{start}$ | $t_{end}$ | Score |
|---|---|---|---|---|
| 24 | 127 | June 10th, 2018 9:43 am | June 10th, 2018 12:50 pm | 0.98 |
| 113 | 202 | June 11th, 2018 7:06 pm | June 12th, 2018 11:18 am | 0.96 |
| ... | ... | ... | ... | ... |
| 721 | 631 | Aug 12th, 2018 1:12 pm | Aug 13th, 2018 4:50 pm | 0.86 |

Machine Learning Pipeline for TS Anomaly Detection

Time segmenting and aggregating → Missing value imputing → Training data extracting with roll windows → Prediction or reconstruction model learning → Discrepancy computing → Anomaly finding
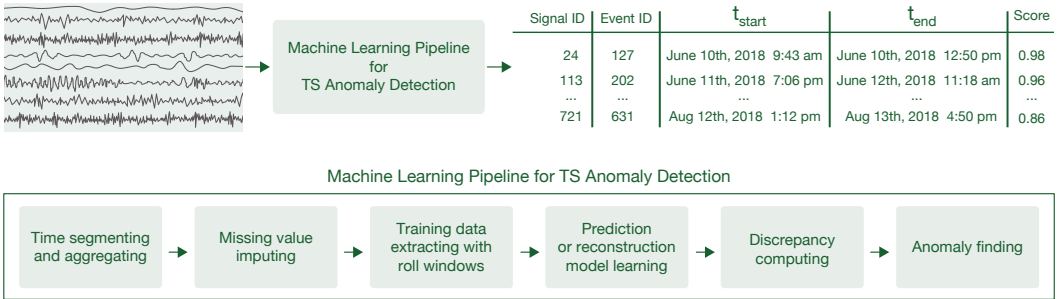
Fig. 3. The end-to-end ML pipeline for multivariate time series anomaly detection. (Top) The pipeline takes multiple signals as inputs and generates a list of intervals of timestamps $(t_s, t_e)$ suspected to be anomalous. (Bottom) The pipeline is composed of six primitives, each serving one particular function.

Fig. 3 describes an entire ML pipeline from start to finish. At a high level, the pipeline takes multiple signals as inputs and generates a list of intervals of timestamps $(t_s, t_e)$ which are suspected to be anomalous (**T1**). The pipeline contains six fundamental primitives [62], each serving a

particular function, and each with various hyperparameters, the settings of which influence the ML outputs[1]. Different settings will result in different ML outputs that are meaningful depending on different analytical demands. Because domain experts generally do not have much ML expertise, we have designed our system around what they hope to observe (**T2**), and provided them with control through the Landing Page (Section 4.4.1).

**Data pre-processing.** The first two primitives transform the data into a clean format $\mathbf{x} = [x^1, x^2, \ldots, x^T]$, where $T$ denotes the total number of time steps and the intervals between $x^{t-1}$ and $x^t$ are equal. To fill any missing values, various predefined imputation strategies are supported, with the mean value used by default. For spacecraft telemetry data, we use intervals from 6 minutes to 6 hours and a zero-order hold strategy[2] to fill missing values according to the domain experts' request. The third primitive prepares a collection of training samples using a sliding window sequence approach, and produces $N$ subsequences $X = \{x_i^{1 \ldots t}\}_{i=1}^N$, where $N = (T - t)/s$ , where $t$ and $s$ represent window size and step size respectively.

**Modeling.** The fourth primitive learns a prediction/reconstruction model to generate a predicted/reconstructed time series $\hat{\mathbf{x}}$. We identify the need to use different algorithms in different scenarios while collaborating with domain experts. MTV has integrated three popular models: the GAN model (TadGAN [21]) performs well at handling complex signal data with myriad fluctuations; LSTM [36] is a well-established sequence learning method suitable for many general scenarios; and Arima [54] is a statistical model proven to work excellently when time series have remarkable trend and periodical patterns. Bearing this high-level knowledge in mind, experts can choose those ML outputs that apply the most suitable model for further investigation (Section 4.4.1).

**Finding anomalies.** Next, the fifth primitive computes the discrepancies between $\mathbf{x}$ and $\hat{\mathbf{x}}$ to locate potential anomalies, under the logic that higher discrepancies suggest a higher chance that the segment is anomalous. In other words, we generate $\mathbf{E} = [e^1, e^2, \ldots, e^T]$ to measure the difference at every time step. Then we apply an exponentially weighted moving average (EWMA) on it, obtaining the smoothed error as below:

$$\mathbf{E_s} = [e_s^1, e_s^2, \ldots, e_s^T] \tag{1}$$

The last primitive takes this error sequence as an input and computes a threshold. Any values regarding the smoothed errors above the threshold are considered to be anomalies. The threshold is selected from the set: $\theta = \mu(\mathbf{e}_s) + k\sigma(\mathbf{e}_s)$ where $\mu$ and $\sigma$ denote the mean and standard deviation respectively. Eventually, $\theta$ is determined by finding a threshold that would bring about the maximum percent decrease in the mean and standard deviation of $\mathbf{e}_s$ if all error values above the threshold were eliminated [36]. Now each anomalous sequence $\mathbf{e}_{seq}$ (continuous sequences of $e_s^i$ with the value above $\theta$) can be assigned a severity score $s$:

$$s = \frac{max(\mathbf{e}_{seq}) - \theta}{\mu(\mathbf{e}_s) + \sigma(\mathbf{e}_s)} \tag{2}$$

Now we can represent each anomalous sequence $\mathbf{e}_{seq}$ in the format of $(t_s, t_e, s)$, where $t_s$ is the starting time of this sequence, $t_e$ is the ending time, and $s$ is the assigned severity score.

## 4.2 Improve Machine Learning with Annotations

An unsupervised ML pipeline detects anomalies that generally show up as unexpected temporal patterns within certain time periods. As shown in Fig. 4, in the early phase when there is a

---

[1]The url (https://bit.ly/mtv_ml_pipelines) links to our github repo (**anonymized for review purpose**) where we detail how the hyperparameters of each primitive of the LSTM pipeline are set. More pipelines can be found under folder `XXXX/pipelines/verified/`.

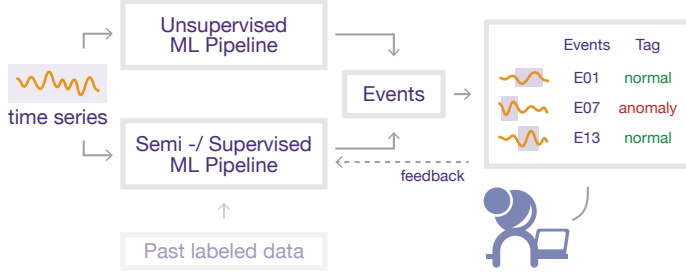[2]Zero-order hold: https://en.wikipedia.org/wiki/Zero-order_hold

Fig. 4. An illustration of how MTV learns from expert feedback and improves ML performance over time.

lack of annotated data, an unsupervised ML pipeline is used to locate anomalous events. These anomalies are presented to experts for annotation. The annotated events are then provided to the semi-/supervised LSTM pipeline, which can be optionally pre-trained with past labeled data. This pipeline keeps learning from feedback and improving on performance over time (**T6**).

We base the refreshment process of the semi-/supervised pipeline on application-specific batch processing of annotations. We justify our design based on the high variability between domains in terms of the frequency of anomalies. For example, based on the satellite company's particular needs, we decided that a weekly update is sufficient.

## 4.3 Enhance Annotating Efficiency with Shape-Matching

Due to the large-scale nature of the data, experts are only able to annotate a limited number of events. We use the idea of "shape-matching" to enhance annotating efficiency (**T5**), as well as to improve the performance of the ML pipeline (**T6**). We introduce a novel shape-matching algorithm for efficient shape search within time series data.

The algorithm is performed in a signal-based manner and outputs a set of candidate shapes in a signal $\mathbf{x}$ that are similar to a particular sub-sequence $\mathbf{s}$, under the constraint that the returned segments will not overlap with each other or any existing events. Such constraint is important without which experts will feel confused by overlapped events. We adopt a sliding window approach to generate subsequences of $\mathbf{x}$, then compare them to $\mathbf{s}$ through a similarity measure $f(\cdot, \cdot)$. We chose $f$ to be the total cost of finding the optimal mapping between two sequences, using either Euclidean or Dynamic Time Warping (DTW [11]) as the similarity measurement. The Euclidean method tends to find exactly matched shapes, while DTW can tolerate a certain level of shifting.

The full algorithmic process is described in the following pseudo-code (algorithm 1). At each checkpoint $t_c$, we attempt to include the current most similar shape $\mathbf{c}$ in the candidate subset. Prior to that, we check whether $\mathbf{c}$ overlaps with a preexisting candidate shape; if so, we keep the most similar of the two. This procedure will return a set of non-overlapping candidate shapes $C$ that are most similar to $\mathbf{s}$.

We argue for three main uses for shape-matching (described below), all of which greatly enhance annotating efficiency and make anomaly investigation and annotations feasible for large-scale data.

**Annotation sharing.** Once an expert confirms that an annotated sequence is indeed an anomaly, s/he will be inclined to search for other segments that are similar to the confirmed event, but have not yet been flagged. These similar segments are likely to receive the same annotations as the confirmed one. Fig. 7ⓑ shows the selected event (b1) and its similar segments (b2 and b3). Experts are able to quickly assign tags to them one after another, or even all at once (Fig. 6ⓒ — "override all segments tags").

---

**Algorithm 1** Shape Matching

---

1: **Input:** Signal $\mathbf{x} = [x^1, x^2, \ldots, x^T]$
2: Sub-Sequence $\mathbf{s} = [s^{k+1}, s^{k+2}, \ldots, s^{k+t}]$ where $1 \leq k < (T - t)$.
3: **Output:** Candidate sub-sequences $C$.
4: Initialize $\mathbf{c}$ = None                                          // to hold candidate shape
5: Initialize $t_c = t$.
6: **for** $i = 1$ **to** $T - t$ **do**
7:     $\tilde{\mathbf{x}} = [x^i, x^{i+1}, \ldots, x^{i+t}]$                                          // subsequence of $\mathbf{x}$
8:     **if** $f(\tilde{\mathbf{x}}, \mathbf{s}) < f(\mathbf{c}, \mathbf{s})$ **then**
9:         $\mathbf{c} \leftarrow \tilde{\mathbf{x}}$                                          // update $\mathbf{c}$ to hold a more similar shape
10:     **end if**
11:     **if** $i > t_c$ **then**
12:         **if** $C \cap \{\mathbf{c}\} \neq \varnothing$ **then**
13:             $\mathbf{v} \leftarrow C \cap \{\mathbf{c}\}$                                          // overlapping sequence
14:             $\mathbf{c} \leftarrow \mathrm{argmin}\,[f(\mathbf{c}, \mathbf{s}), f(\mathbf{v}, \mathbf{s})]$                                          // keep most similar
15:         **end if**
16:         $C = C \cup \{c\}$                                          // add candidate shape
17:         $\mathbf{c}$ = None                                          // reset
18:         $t_c = t_c + t$
19:     **end if**
20: **end for**
21: **return** $C$

---

**Decision support.** Sometimes the status of an anomaly is ambiguous, and choosing a proper tag is difficult. In this case, an expert can use shape-matching to find similar shapes (non-anomalous segments) to this event. By comparing the differences between them, experts can quickly clarify why the anomaly detection algorithm identified this segment as anomalous. We introduce a superposed visualization in Section 4.4.7 that allows for such comparisons. Fig. 6ⓒ shows one example, where the event under investigation can be compared to similar segments.

**False alarm mitigation.** After events have been properly annotated, we are able to mitigate false alarms (false positives) by leveraging the shape-matching algorithm. Consider an event that experts may think is unworthy of exploration — say, one tagged as "Do not investigate" or "Normal" (see details in Section. 4.4.6). We can ask the system to prune existing events similar to this one using the method described below:

We use the shape of one such event as a template to search a set of candidate similar shapes, which can be tuned according to a pre-defined threshold. Those supposedly anomalous segments identified by the ML pipeline are now pruned by checking whether they overlap with these candidate similar shapes, thus reducing the number of false positives.

## 4.4 Visualization and Interaction

Our interface consists of two pages: the Landing Page and the Investigating Page. An expert starts on the Landing Page (Fig. 5), which provides an overview of anomaly detection results and allows experts to select certain signals and ML results for further investigation. The Investigating Page (Fig. 1), MTV's most important component, contains three major panels: the Signal Overview ⓐ, the Signal Focused View ⓑ, and the Side Panel ⓒ, with four collapsible sub-views included. All three panel views work in a synchronized fashion to support the workflow described in Table 1.

*4.4.1* ***Landing Page****.* Given a time series dataset consisting of many signals, a machine learning expert will explore a collection of ML pipelines and hyperparameter settings to identify anomalies. The results of these pipelines are then stored in the database (**T1**). However, an individual or team of experts may be interested in monitoring only a subset of signals and results. This selection is often driven by the experts' particular expertise and the problems they're interested in tackling. The Landing Page is designed to offer them this flexibility (**T2**).
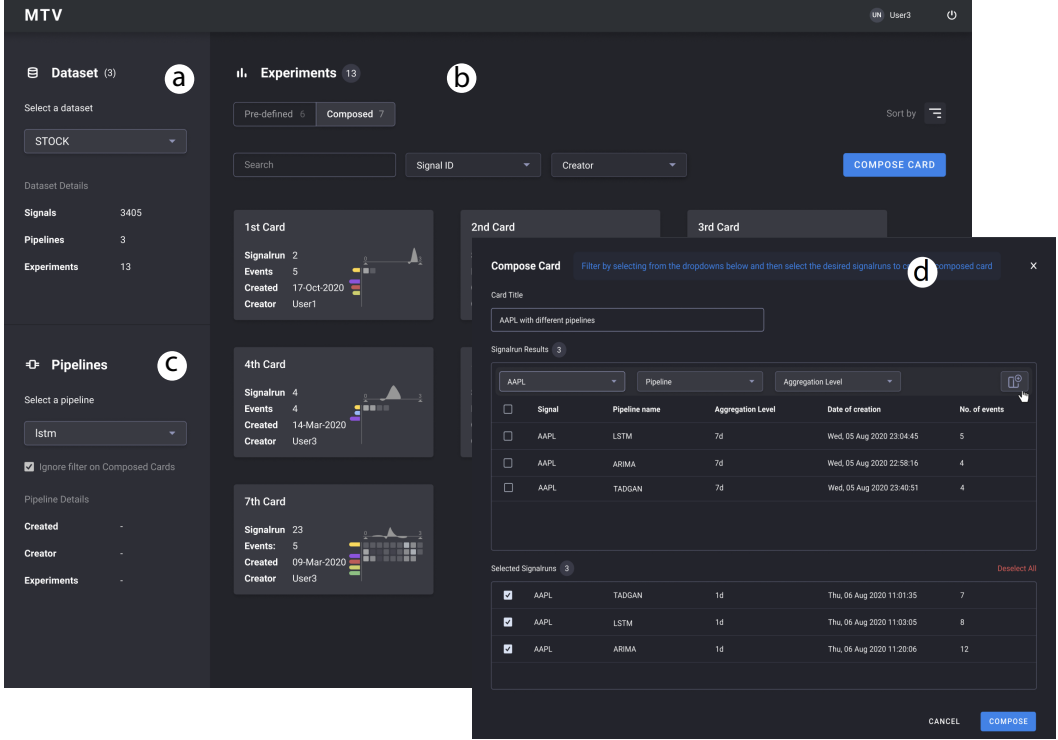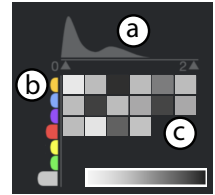


Fig. 5. The Landing Page. The Dataset Panel (a) allows users to select a dataset and observe its basic information. The Main Dashboard (b) lists all the experiment cards following a user-determined order and can be filtered by whether certain pipelines have been applied in the Pipeline Panel (c). Users can click the "COMPOSE CARD" button to enter the Compose Card Window (d) where they can select ML outputs of interest to create a new experiment card for future collaborative analysis.

An expert can choose different signal datasets to explore. In Fig. 5ⓐ, we see that the STOCK dataset contains 3405 signals (i.e., stocks), 3 pipelines (i.e., ARIAM, TADGAN, and LSTM) and 13 experiments. An experiment is defined as the outputs of one ML pipeline on a certain signal subset.

We propose a novel card-style visualization to display the summary information for each experiment. The "experiment card" shows three valuable pieces of information at a glance (shown on the right). The area chart ⓐ at the top encodes the distribution of severity scores according to (Eq. 2). The colored bar chart ⓑ on the left side shows the number of events and their corresponding tags (Sec. 4.4.6). By checking the information encoded in the previous two charts, it is possible to obtain quick insights about the overall



severity and urgency of this analysis task and track its progress. The matrix ⓒ offers experts

information about how anomalies are distributed over signals. Each cell of the matrix represents one signal. The cell color encodes the number of anomalies detected in that signal, and a darker shade indicates a larger number of detected anomalies.

All experiment cards are listed in the right area of the Landing Page (Fig. 5ⓑ) following a user-determined order. Various filter types, sorting by signal IDs, creator and/or card title, can be applied to facilitate efficient card exploration. In particular, the cards can be filtered by whether certain pipelines have been applied (Fig. 5ⓒ).

To provide further flexibility, the Landing Page also allows experts to select their ML pipeline and hyperparameter settings by exposing a minimal set of options to compose a new experiment card. These options are intuitive and easily understandable. They include several important parameters such as the ability to select the level of aggregation, the strategy for imputing a missing value, or a particular ML algorithm. Fig. 5ⓓ shows an example experiment card, created by an expert and titled "AAPL with different pipelines." This card is meant to compare the difference in results when running three different pipelines on AAPL (Apple) stock, with an aggregation level of one day.

Finally, we stress the importance of the Landing Page, as it is an essential step (**T2**) of the workflow for general time series anomaly analysis. However, the design of this page could be domain-specific. Fig. 5 shows a prototype version that we implemented for experimental purposes. However, for different application scenarios, the names of UI components in Fig. 5 could be changed, and the set of filtering options in ⓓ varied.

*4.4.2*   ***Signal Overview****.* Once experts select their "experiment card", they go to the Investigating Page for further analysis (Fig. 1). The first view in the Investigating Page is the Signal Overview. This view (Fig. 1ⓐ) presents experts with an efficient way to scan the dynamics of every signal, as well as how anomalies are distributed. This helps experts make informed decisions about which signals to investigate further (**T3**).

Given limited screen space, the design of multivariate time series should be highly space-efficient. Therefore, we choose *small multiples*, a space-efficient technique that is widely used for visualizing multiple time series [42]. The timeline of each signal is aligned horizontally. In addition, we opt to use a line chart rather than another chart type (e.g., area chart, horizon graph), because of its interpretability. Line charts are well-suited for point-wise inspection. Additionally, line charts make trend lines very apparent, making it possible to visualize anomalies through simple trend tracking [42]. We visualize anomalous events by highlighting the curves with a warning color (Fig. 1-a1). Events across signals that appear close together in time, known as co-occurring patterns, may indicate a larger event. This further suggests that the signal order has a significant impact on experts' ability to investigate. Thus, we propose a novel layout algorithm to optimize the order.

**Order optimization:** Our goal is to put "similar" signals as close together as possible. Similarity here refers to the length of overlapping anomalous events present in two signals. To that end, we present a novel metric to measure the similarity between two signals $\mathbf{x}_a$ and $\mathbf{x}_b$:

$$Sim_{\mathbf{x}_a,\mathbf{x}_b} = \frac{|\mathbf{A}_a \cap \mathbf{A}_b|}{\min(|\mathbf{A}_a|, |\mathbf{A}_b|)} \tag{3}$$

where $\mathbf{A}_a$ and $\mathbf{A}_b$ are the sets of anomalous events detected in the first and second signal respectively, generally represented as $\mathbf{A} = \{(t_s, t_e)^i \mid i = 1, \ldots, K\}$ with $K$ denoting the number of identified anomalous events in a particular signal. In addition, $|\mathbf{A}| = \sum_{k=1}^{K}(\hat{t}_e - \hat{t}_s)$ denotes the total length of the events. The intersection of the two event sets, denoted by $\mathbf{A}_a \cap \mathbf{A}_b$, is the event set containing all timestamps from $\mathbf{A}_a$ that also belong to $\mathbf{A}_b$. For example, when $\mathbf{A}_a = \{(1, 5), (8, 12)\}$ and $\mathbf{A}_b = \{(4, 6), (9, 11)\}$, their intersection should be $\{(4, 5), (9, 11)\}$ whose length is $1 + 2 = 3$.

This metric is inspired by the *Overlap Coefficient* [66], which measures the overlap between two finite sets, defined by the equation $overlap(X, Y) = \frac{|X \cap Y|}{\min(|X|,|Y|)}$. It is known that the Jaccard Similarity becomes inefficient when two sets vary significantly in size, but the Overlap Coefficient overcomes this issue. If set $X$ is a subset of $Y$, or the converse, the overlap coefficient is equal to 1. Our defined metric $Sim_{x_a,x_b}$ shares the same advantages.
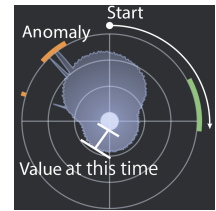
After defining the similarity metric between signals, we can now use dimensionality reduction techniques to put the signals in sequence. In our case, we choose t-distributed stochastic neighbor embedding (T-SNE [48]) to represent all signals as single-value embedded codes. The signals are then sorted from top to bottom in the Signal Overview by their code values.

*4.4.3* **Signal Focus View**. This view (Fig. 1ⓑ) extends a piece of the timeline from the selected signal in the Signal Overview, displaying more information (**T4**). Experts can either mouse over context charts (Fig. 1-a2) or use the zoom and pan functions in the Signal Focus View for flexible exploration. As shown in Fig. 1-b2, anomalies are highlighted such that the color of the line corresponds with the color of the tag associated with the event. To further enhance awareness, the color of the header bar double-encodes this tag information. In addition, a transparent grey background is added to make anomalies more visually apparent. Prediction results are visualized with a thinner curve in bright yellow.

The smoothed errors (Eq. 1) are represented as a centered flow on the top of the chart, where a thicker flow width indicates larger discrepancies between the original values and the predicted values (Fig. 1-b1). This visualization increases the transparency of the model and enables experts to visually evaluate the quality of the model as well as interpret why a certain anomaly was identified by the ML pipeline. Error flows can also be used by experts as a visual clue, guiding them to interact with (CRUD) anomalies (**T5**).

*4.4.4* **Periodical View**. This view (Fig. 1ⓒ) is designed for analyzing periodic patterns of the selected signal (**T4**). The table on the top (Fig. 1-c1) summarizes the number of overall tags, as well as tags per year and per month, for a signal in focus. The bottom graph (Fig. 1-c2) provides experts with a new perspective for exploring the periodical patterns of a signal. Three levels of periodical glyphs, corresponding with the three different time granularities (i.e., year, month, and day), are proposed to support multi-scale analysis.

The glyph design (shown on the right) is inspired by the circular silhouette graph [2]; the glyph employs a polar coordinates system, where the angle encodes the time point in a year/month/day and the radius indicates the value at this time point. We propose the use of *small multiples* [42] due to its spatial efficiency and the ease of side-by-side comparison. If for one time period (year/month/day) a signal has an irregular shape or unusual spikes, this may indicate anomalies. In addition, we highlight the anomalies, according to their tags, by overlaying radial segments in the corresponding time periods. This enables experts to observe how anomalies are distributed across years/months/days. For example, in Fig. 1-c2, the green-tagged event and red-tagged event occurred near to the start or the end of the year, while the three orange-tagged events are close to the middle of the year.

*4.4.5* **Signal Annotations View**. This view (Fig. 6ⓐ) provides an overview of all tags across the currently selected signal. From here, experts can quickly glance through the event information (starting and ending times) and what tags are associated with these events (**T4**). We follow the design of GitHub issue labels to visualize the tag — a rounded rectangle with the background encoding the tag type and the text showing its associated name and meaning. Experts can click to open one event and explore the most recent annotations associated with the event. To improve

efficiency, experts are allowed to directly post their comments or assign a tag here (**T5**). These events are chronologically ordered (from top to bottom in Fig. 6ⓐ), and the sequence corresponds with the order (from left to right in Fig. 1ⓑ) on the focused view, that is, green-orange-orange-orange-red.

*4.4.6*   ***Event Details View and Multi-Aggregation Viewer***. Experts can either go into the Event Details View (Fig. 6ⓑ) by directly clicking it on the Side Panel (Fig. 1ⓒ) or by using the "go to Event Details" button (Fig. 6ⓐ) from the Signal Annotations View. When the view is opened, the focal chart will be equipped with the Multi-Aggregation Viewer (Fig. 7ⓒ), where experts can choose different granularities to explore contextual information about the current anomaly (**T4**).

The Event Details View, from top to bottom, displays the starting and ending times, tag information, severity score (valid only when the source is "ML"), source, and the comment box. The source can be either "ML", "USER", or "Shape-matching." In this view, the comment box shows all the historical annotations of the event (**T4**), in contrast to the Signal Annotations View, where only the five most recent annotations are shown. Along with the other coordinated views, this view allows experts to perform in-situ annotation and communication (**T5**), with all the necessary contextual information displayed on one screen.

We have designed six general types of tags, plus the status "untagged," to assist collaborations between experts. Here are these tags and their meanings: (1) `Do not investigate` (action tag): "We are not interested in this and have decided not to investigate." (2) `Postpone` (action tag): "This event is interesting but of low priority; we will postpone its investigation." (3) `Investigate` (action tag): "This event is interesting and we should investigate it now." (4) `Problem` (info tag): "This is a new problem, and while we can describe it colloquially, we don't have a term for it yet." (5) `Previously seen` (info tag): "This is a well-known problem that we have investigated before." (6) `Normal` (info tag): "This event is normal, has an obvious explanation, and is not harmful."

The tag design is the result of many design meetings with our domain experts (i.e., P1-P6 and E1-E3). The first three action tags are meant to suggest the next step that should be taken pertaining to an event, while the last three explain what has already been decided. The tag of one event can be changed over time. In practice, if an event is tricky, experts often use action tags initially to facilitate team communication, and switch to a specific info tag later on when they reach a consensus.

*4.4.7*   ***Similar Segments View***. This view (Fig. 6ⓒ) provides experts with the ability to search the most similar segments for a selected event. Assume a number of similar segments (up to 100 by default) are returned by the shape-matching algorithm. The bar chart at the top of this view is used to filter the segments based on the similarity score. Below the chart is a list of segments showing more detailed information, such as start and end time. The graph on the right plots the returned segment line overlaid by the original line, color-coded by its tag (in this case orange, which is consistent with its assigned tag — Fig. 7-b1), allowing experts to visually compare how similar they are. Meanwhile, the corresponding similar segments on the focal charts are highlighted using the dashed white border (Fig. 7-b2 and -b3). As described in Section 4.3, the Similar Segments View can be used for annotation sharing, decision support, and false alarm mitigation, in order to boost annotation efficiency (**T5**).

## 5   EVALUATION

To understand the usability and usefulness of MTV , we carried out two user studies based around two potential stakeholders: a domain expert, who has area expertise and analyzes anomalies as part of his or her job, and a general end-user, who does not have area expertise or anomaly analysis duties, but may still be interested in performing such tasks.
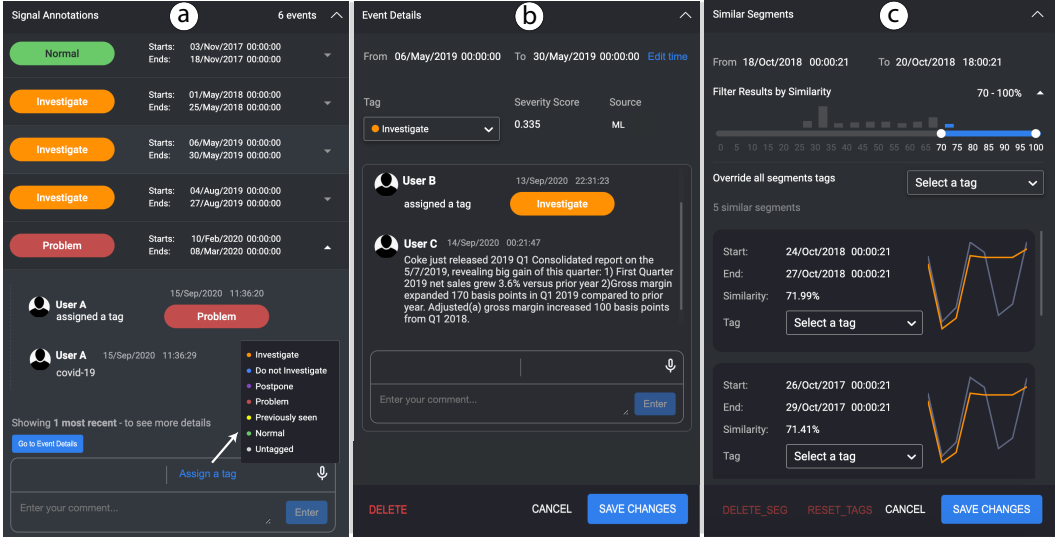
Fig. 6. Three sub-views in the Side Panel: (a) the Signal Annotation View provides an overview of the annotations (ordered by event time) made for the selected signal; (b) the Event Details View shows more details about one particular event, such as severity score and source; (c) the Similar Segments View displays the search results of the shape-matching algorithm and allows users to perform quick annotations.

## 5.1 User Study Using Spacecraft Telemetry Data with Domain Experts

We deployed MTV with real spacecraft telemetry data in a production environment. We carried out four case studies with our expert collaborators (P1-P6) and conducted semi-structured interviews to collect qualitative feedback about MTV.

*5.1.1* ***Experimental Setup****.* We prepared the experiments using 55 real signals (measuring device electrical power, thermal temperature or attitude) tracked over a 5 year period. We ran ML pipelines with a variety of possible settings by changing different hyperparameters that the experts wanted to control, such as the aggregation level (from 6-minute to 6-hour), imputation strategy (mean values or zero-order hold), modeling algorithm, etc. Then we asked the experts to work together to compose four experiment cards using the Landing Page. The team would later collaborate to analyze the four experiment cards.

For each experiment card, the experts selected between 4 and 12 signals to analyze. For the first card, Case 1, four signals total were chosen from a variety of spacecraft subsystems. For Case 2, eight temperature signals were selected together with four environmental signals (e.g., sun elevation). For each of Cases 3 and 4, four attitude-control signals plus two environmental signals, and four electrical-power signals and two environmental signals were chosen.

The experts were then asked to enter MTV within the next 24 hours and analyze these experiment cards. During this period, they can enter MTV anytime to create comments, check comments from other team members, and if necessary, have in-situ discussions under events of interest. Experts were asked to record all on-screen activities while performing these tasks, and to use the think-aloud protocol. After the 24-hour period was complete, we conducted semi-structured interviews with all the experts one by one to collect qualitative feedback.

*5.1.2* ***Results****.* We summarized event and annotation numbers for the four case studies in Table 2. We observed that the experts were able to used their domain knowledge to mark their own suspicious
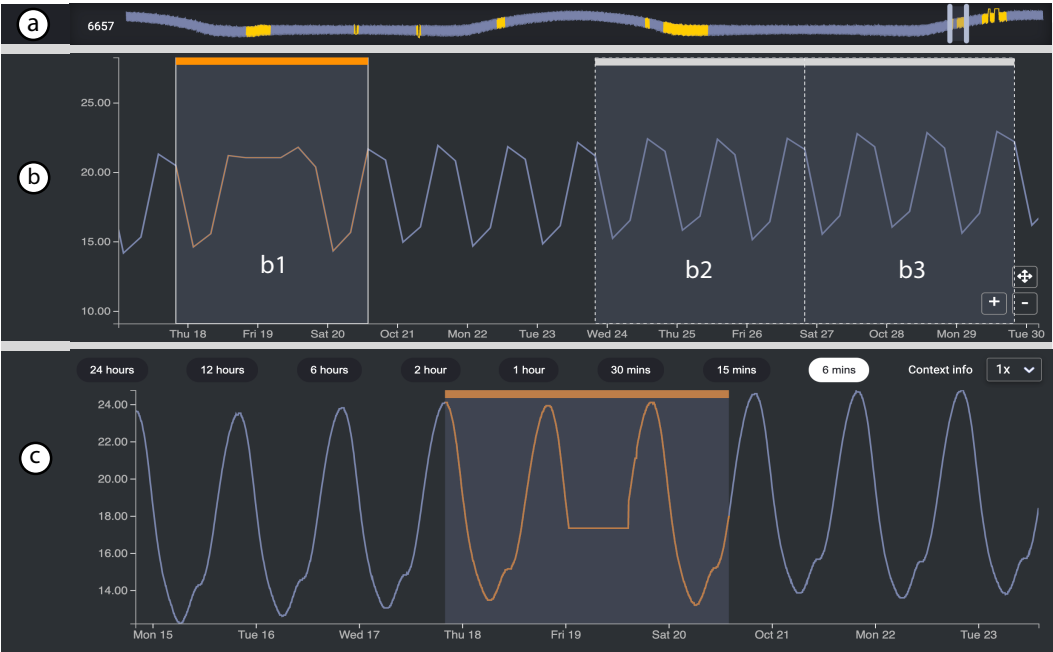
Fig. 7. (a) The overall time series for signal X — a temperature signal from a spacecraft; (b) The interval b1 is identified by ML and tagged as "investigate" by one expert (P2). The expert uses the Similar Segments View (Fig. 6-right) to find similar segments for comparative analysis. (c) The interval b1 is observed at a 6-minute aggregation level using the Multi-Aggregation Viewer. Seeing this, P3 suspects that this anomaly was the result of a time gap issue (missing values at a certain time period). The missing values were then filled in using the last valid value, leading to a flat line in the 6-minute level; this makes shape b1 (in the 6-hour level) stand out, showing an unusual pattern even compared with its most similar segments (b2, b3).
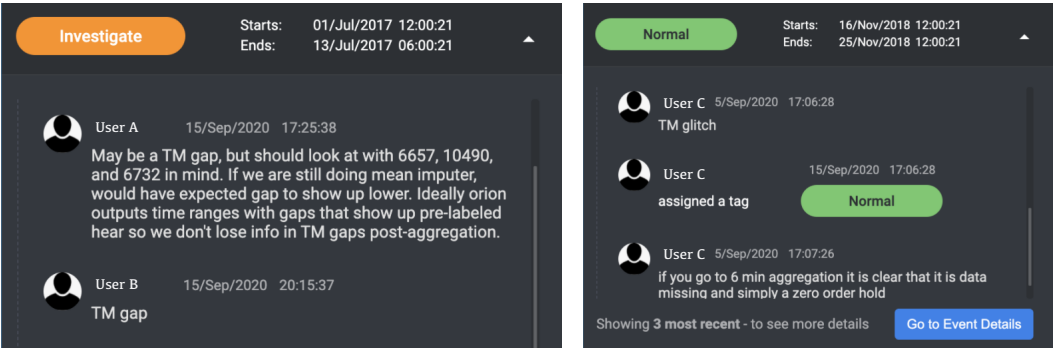


Fig. 8. Example comments demonstrating how satellite experts discussed anomalous events. (Left) User A presented his suspicion that a certain anomaly may have come from a TM gap (missing values). Several hours later, User B (a more senior engineer) confirmed this suspicion and changed the tag from "investigate" to "normal." (Right) User C, who felt confident about his finding, directly tagged the event as "normal" and added an explanation.

| Case | ML Event | User-created Event | Comment (avg.) | Tag (avg.) |
|------|----------|--------------------|----------------|------------|
| 1 | 38 | 15 | 162 (3.1) | 58 (1.1) |
| 2 | 45 | 12 | 87 (1.5) | 60 (1.1) |
| 3 | 40 | 8 | 96 (2.0) | 48 (1.0) |
| 4 | 23 | 10 | 66 (2.0) | 40 (1.2) |

Table 2. Statistics displaying the event numbers and collected annotations for all four case studies. (avg.) indicates the average annotation number per event.

events in addition to those identified by ML. Over the four cases they created 15, 12, 8, and 10 events, respectively.

A large number of comments about one event suggests an active discussion. We observed that the average comment numbers for the four case studies were 3.1, 1.5, 2.0, and 2.0 respectively, indicating that experts can use the system to efficiently perform annotations.

Experts can perform collaborative analysis effectively as well. Fig. 7 demonstrates the high quality of annotations through an example in which P2 investigates a temperature signal X. P2 selected a region in signal X (Fig. 7ⓐ) and tagged the event (b1) as "investigate." He used the Similar Segments View to find similar segments (b2 and b3) for comparative analysis and obtained an initial suspicion: "*Maybe a time gap issue. We should look at its correlated signal Y and Z to double confirm.*" After this comment, another expert (P3) responded, "*Signal Y and Z are fine. I think it should be a time gap issue. If go to 6-minute aggregation (Fig. 7ⓒ), you will find the flat line which actually is caused by missing value filling strategy*". Finally, P4 concluded, "*Missing data issue confirmed. I will switch the current tag from investigate to normal.*" The discussion above demonstrates that experts are able to use MTV to conduct collaborative analysis. Fig. 8 shows two additional examples that demonstrate two different types of collaborating patterns.

*5.1.3 Feedback Summary.* In general, the experts greatly appreciated the system and valued its potential to enhance their efficiency for time series anomaly analysis. They were happy to introduce the current version of MTV to their colleagues for further use and testing. We summarize their feedback and suggestions as follows.

*Knowledge*. P1 and P2 highlighted, "*MTV provides a comprehensive view of telemetry data. It creates a well-defined knowledge base that serves as a reference for the entire team during their investigation process.*" Both are interested in learning how to better organize and leverage existing annotations.

*Scalability*. All experts agreed that MTV makes time series anomaly analysis possible for large-scale time series data. P4 commented, "*MTV saves me a huge amount of effort on finding suspicious anomalies. Without the system, we are only able to monitor few numbers of signals.*" P6 added, "*MTV offers us a fantastic way to share and communicate what we found, before which we wasted too much time in using CSV to do the sharing.*"

*Usability*. All the experts agreed that functionalities provided in MTV allow for efficient event exploration and effective decision making. P1 suggested, "*The error river chart is useful for me to get a sense of how serious a certain event is predicted to be. But there will be a case when all the other events are minor in comparison and cannot be seen. I suggest allowing a zoom function along y axis of this chart.*". P6 commented, "*The Similar Segments View opens me to a brand-new way to investigate anomalies.*"

*Confidence*. MTV facilitates discussions and allows the team to share insights and conduct collaborative analysis. P3 was excited about this feature, "*When I can document and share what I thought in such an organized way, I would gain more confidence in my annotations.*"

## 5.2   User Study Using Stock Data with General End-Users

To obtain a more comprehensive understanding of how general end-users perceive MTV, we conducted experiments using stock price data. We selected 10 stocks from different sectors, such as energy, technology and finance, and used their daily price data since 2015. Hence, each signal (i.e. stock) contains around $2,000$ data points.

*5.2.1*   **Participants.**  We recruited 25 participants (18 male, 7 female; aged 23-40) via email invitations and on-campus advertising at 3 universities. We did not set many constraints when selecting participants. Each participant had between 0 and 15 years of data analytics experience ($\mu$=4.52, $\sigma$=4.59) and between 0 and 10 years of machine learning experience ($\mu$=2.42, $\sigma$=2.46). They also had varying occupations, from students and consultants to data analysts and UI/UX designers. All showed a strong interest in analyzing time series data from their daily lives. Only eight participants out of 25 had experience in the stock market. Because real-world users of MTV will likely have similarly diverse backgrounds, we wanted to see whether our volunteers found the system both usable and useful.

*5.2.2*   **Experimental Setup.**  We sent each participant a website link to perform the user study. The website included detailed instructions to guide them through the experiment. The participants started the studies asynchronously in random order. The entire process was meant to last 1 hour. At the start, we asked the participants to fill out a background information questionnaire. The study itself began with a 10-minute training, after which participants were asked to perform three exploration tasks and one case study. During the training, each participant watched a tutorial video, performed three exploration tasks in MTV , and answered questions about each task to confirm their understanding of the core concepts and features. During the tasks, participants were asked to follow step-by-step instructions in order to explore data using particular system features. After each task, they were asked to answer several questions about what they had found (such as the number of anomalies in stock A), as well as to provide feedback on whether a certain feature was easy to understand and use.

After completing the three exploration tasks, the participants were asked to collaboratively work on one open-ended task — a case study — creating annotations and adding their own interpretations in an asynchronous manner. The case study involved investigating three stocks (COKE, REGN, and INTC) from different sectors — consumer, healthcare, and technology — together with the Nasdaq Composite (Fig. 1ⓐ). The participants were able to access stock news websites or use search engines to help with the annotation. They could choose to annotate any number of events of interest, with no minimum. Finally, they filled in post-study questionnaires regarding the effectiveness and the usability of our system, and went through a short (5 to 10 minute) informal interview in the form of one-on-one Zoom meeting.

*5.2.3*   **Results**.  Our aim with this case study was to understand whether participants could create good annotations and make sense of ML results using MTV. Our participants took 12.7 minutes ($\sigma = 7.2$) on average to finish their annotations during the case study. 27 events were identified by the ML algorithm, of which 25 (92.6%) events were annotated (either tagged or commented). 11 additional events were manually created by users. In total, the participants created 65 tags (1.7 per event) and 128 comments (3.4 per event). We found that while a few participants did not create any annotations, this seeming lack of engagement was actually collaborative — these users explained that they saw other people's annotations and thought they were reasonable, and felt no need to give extra explanations.

To evaluate the quality of annotations, we sought help from 2 external volunteers with more than 3 years of stock market experience. Working together, we marked any tags or comments
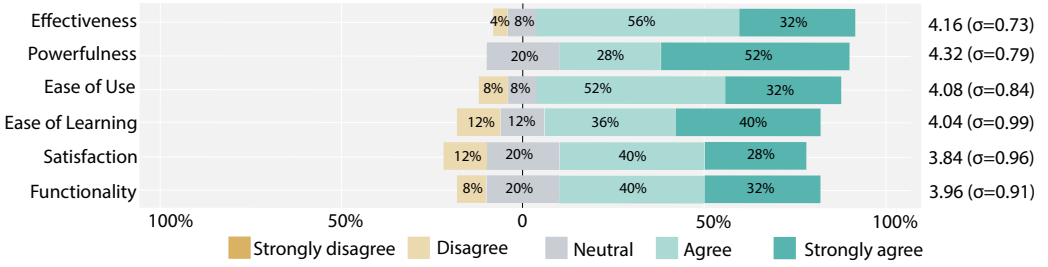
Fig. 9. Ratings of the overall MTV experience by general end-users.

that did not make sense as "invalid." After this, 93.8% (61/65) tags and 91.4% (117/128) of the total comments were considered valid. As an example, Fig. 6(a) and 6(b) show real annotations made by users, which correspond to the two anomalies highlighted in Fig. 1-a1 (REGN) and -b2 (COKE) respectively. Interestingly, Fig. 1-b2 was first tagged by User B and later commented on by User C. The comment "*the rise is because that COKE just released 2019 Q1 Consolidated report on the May 7th 2019...*" clearly explains the reason for this abnormal rise in the stock price of COKE.

In conclusion, we confirm that participants can use MTV to create good annotations to make sense of, or even improve upon, ML results.

**User experience — overall assessment**. In the post-study questionnaires, we used 5-point Likert-scale questions (1 = strongly disagree and 5 = strongly agree) to assess the overall functionality and usability of MTV. Fig. 9 shows their overall ratings. We found that in general, participants highly rated their experience with MTV .

Notably, ratings regarding usefulness are near the top of the scale, with the highest score ($\mu = 4.16$, $\sigma = 0.73$) for powerfulness (in annotating and sharing annotations), and the second-highest score ($\mu = 4.16$, $\sigma = 0.73$) for effectiveness (at detecting and investigating anomalies). This suggests that MTV has well achieved its goals of supporting time series anomaly detection, investigation, and collaborative annotation. One participant (User B in Fig. 6(b)) commented, "*MTV is amazing tool. I love the collaboration aspect very much. I will be more confident to make my annotation when I can check other people's comments. And I felt excited when my annotation is supported by other people's comments.*". Only one person (4%, 1/25) thought our system was not effective enough to support his investigation, because he was not able to access more fine-grained price data (e.g., hourly). We think this is a flaw of the data and not directly related to our system.

The ratings for the usability aspects were slightly lower than the usefulness scores, but still good. The participants rated ($\mu = 4.08$, $\sigma = 0.84$) for ease of use, ($\mu = 4.04$, $\sigma = 0.99$) for ease of learning, and ($\mu = 3.84$, $\sigma = 0.96$) for satisfaction (in using the system) which indicates room for further improvement. Some participants complained about the difficulty of learning all the interaction logic for such a comprehensive system.

Users rated the functionality (a.k.a. utility) at ($\mu = 3.96$, $\sigma = 0.91$). A few participants (8%, 2) thought the features provided by MTV did not cover all their needs. For example, one wrote, "*I would like to have an option to refer to another event when adding a comment to a new one.*". Another one said (refer to Fig. 1-a1), "*Looking at stock Regeneron in 2020, it has many anomalies that are nearby each other. It makes sense to merge them into one because they are caused by the Covid-19 vaccine race. So I want a function to support nearby anomaly merging*".

**User experience — features ratings**. We also asked the participants to rate how helpful each feature was to the investigation and annotation process during the case study. The ratings are listed in Fig. 10. Overall, the scores hovered around 4 — "very helpful" — which demonstrates users' contentment with the features.
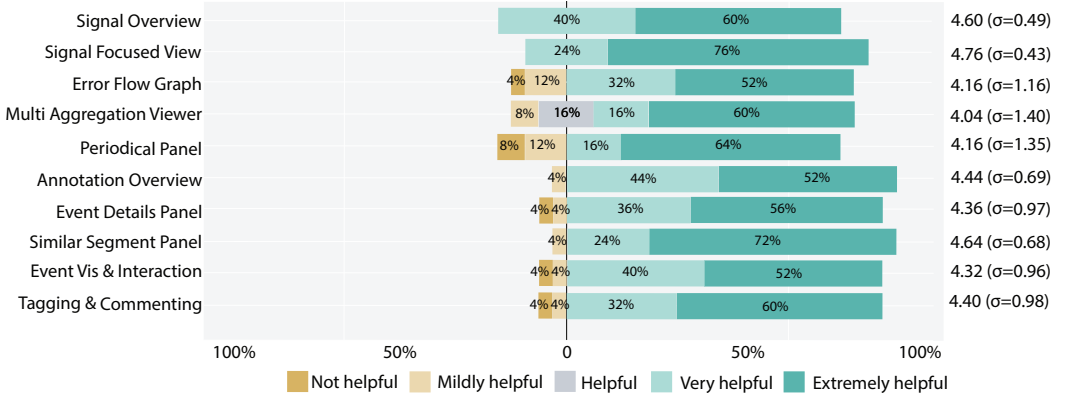
Fig. 10. Ratings of individual MTV features by general end-users.

The Signal Overview ($\mu = 4.60$, $\sigma = 0.49$) and the Focused View ($\mu = 4.76$, $\sigma = 0.43$), the two most important and frequently used features, were rated near the top of the scale. All participants thought these two views were very or extremely helpful. Nearly all participants (96%, 24/25) agreed or strongly agreed on the usefulness of the Similar Segment Panel. One commented, "*The function of similar shape search is so novel to me and this panel is so well designed. It really helps me verify or spread my annotations efficiently!*" Another interesting observation is that people's opinions on the Periodical Panel ($\mu = 4.04$, $\sigma = 1.40$) are slightly polarized, showing that insights regarding periodical patterns are valued differently by different people.

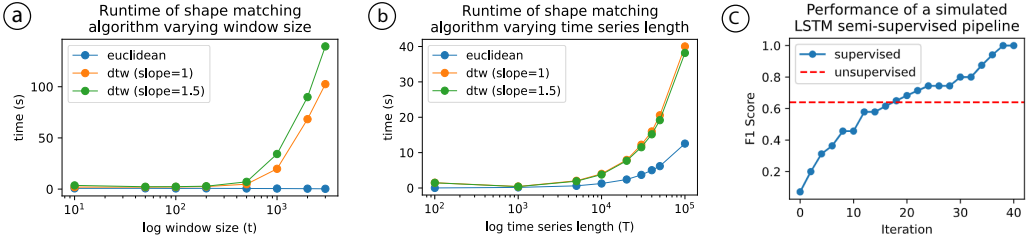## 5.3   Quantitative Evaluation of Algorithms



Fig. 11. (a) Time performance of shape-matching with fixed time series length 5K and varying window size (log scale); (b) Time performance of shape-matching with fixed window size 10 and varying time series length (log scale); (c) Performance of our semi-/supervised pipeline (Fig. 4) using a simulated annotation procedure.

**Time performance of shape-matching.** In theory, the time complexity of Algorithm 1 is $O(TD)$, where $D$ is the time needed to calculate the distance between two sub-sequences of length $t$ (i.e., the window size of the queried shape), and $T$ denotes the total length of the time series. In the case of Euclidean distance $D$ is linear ($D = t$), while in the case of vanilla DTW $D$ is quadratic ($D = t^2$). We consider enabling faster variations of DTW, such as the Sakoe Chiba band and the Itakura parallelogram, which make the time complexity dependant on the band width and slope of the parallelogram respectively [41, 56]. In real-world scenarios, we expect $t \ll T$; thus, the algorithm runs in close to linear time.

Fig. 11ⓐ and 11ⓑ report our experimental results on synthetic data when different distance metrics are applied. We generated signals of variable length $T$ of sine waves with added noise sampled from a Gaussian distribution $\mathcal{N}(0, 1)$. Window size denotes the number of data points of the queried shape ($t$), and time series length indicates the total length of the signal ($T$). Note that the x-axis uses a log scale. From Fig. 11ⓐ, we observed that the algorithm would generally take less than 2s when the window size is smaller than 300 and the time series length is fixed at 5K. From Fig. 11ⓑ, we noted that the time cost is less than 5s when the length of the time series is less than 10K and the window size is fixed at 10. In practice, we assume the anomalous window will be small (a few to dozens of data points). Hence, our algorithm is acceptable for real-time interactions in most scenarios.

**Human feedback evaluation.** To validate the function of integrating human feedback back into the pipeline to improve ML performance (**T6**), we conducted a simulated experiment. For simplicity, we assumed that experts are able to annotate $k = 2$ events in a single iteration and tag an event either as "unwanted" (indicating that it should not be identified by ML) or "wanted" (indicating that it should). The simulation stopped when all events had been annotated.

A semi-supervised LSTM pipeline was trained on sequences that had already been labeled by annotators as either anomalous or normal. We used a 70/30 data split on the NAB dataset[3] for training and testing. The training data encompasses 70 events, while the test data has 32 events. Results are shown in Figure 11ⓒ where we observed that the F1 score[4] of a semi-supervised pipeline surpasses that of an unsupervised pipeline when sufficient annotations have been obtained. In addition, observing several flat segments in the figure, we noted that some annotations may not contribute to improving detection.

## 6 DISCUSSION AND IMPLICATIONS

Feedback from our user studies led us to a set of design implications. We divide our discussions into three topics: visualization and interaction, machine learning, and applicability.

### 6.1 Visualization and Interaction

**Visual scalability of the Signal Overview.** Very occasionally, experts may compose an experiment card with many signals. Given limited screen space, the function of the Signal Overview may then be suppressed, as experts can only observe a limited number of signals at once. Although our proposed order optimization method can keep the "similar" signals close to enable better observation of co-occurring patterns, the analysis of large number of signals simultaneously can still be challenging. A more scalable and intuitive visual design should be considered — one with a more information-dense design like LiveRAC [50], or more advanced interactions like `focus+context` [43]. It is worth noting that we should be careful not to increase the complexity of the original streamlined investigation workflow.

**Alternatives for visualizing multiple time series.** Providing sufficient contextual information is crucial for anomaly annotation. However, how users perceive such information depends heavily on how the time series is visualized. Although the line chart is our current choice for visualizing multiple time series, other visualizations, such as area charts, horizon charts, and color-fields, may also be suitable in different contexts [22, 42]. For example, we found that the satellite experts preferred the step chart — a variant of the line chart — because it made it easier for them to track time gaps when a zero-order hold strategy was employed. In addition, a few participants

---

[3]The Numenta Anomaly Benchmark (NAB) dataset is a well-known public time series dataset that contains 45 signals with 94 labeled anomalies. https://github.com/numenta/NAB

[4]F1 score is a measure of a model's accuracy on a dataset which ranges from 0 (the worst) to 1 (the best).

in the stock experiments found it difficult to perceive nuanced differences between multiple time series. To avoid this situation, colorfields or horizon graphs might be considered as alternatives, depending on whether temporal warping or amplitude change are tolerable [22]. As future work, we will explore the best way to integrate other types of multiple time series visualization into our exploration workflow.

**Flexible tag creation.** Our current system does not support flexible tag creation. One obvious benefit of such a function is that experts could use a word or phrase to mark the reason for an anomaly (e.g., lunar eclipse). However, flexible tagging also introduces additional challenges for tag management as each person may have his/her own tagging "language". One follow-up challenge is figuring out how to effectively utilize these tags to improve machine learning models and ensure mutually beneficial cooperation between humans and AI.

**Anomaly comparison across signals.** The current system does not support direct comparison of anomalies from different signals. Because experts must switch between different signals in the Signal Overview, it is also difficult for them to manually perform such comparisons. One easy solution may be the addition of a pattern screenshot function, so that when experts find useful patterns, they can save them as images for later image-level comparative analysis — although the interactivity of this design may get lost as well. This would be an interesting direction to explore.

## 6.2 Machine Learning

**Learning from annotations.** Currently, only info tags are fed back to the ML pipeline to enhance prediction performance. In Section 5.3, we used a simulated experiment where we assumed an expert can only mark whether s/he wants a particular event to be identified in the future. In MTV's tagging system, we treat Normal as equivalent to "unwanted," and both Previously seen and Problem as "wanted." The other three action tags (Do not investigate, Postpone, and Investigate) are designed to facilitate the collaborative progress of teams, and are not used in the feedback loop. The semi-/supervised ML pipeline can also support multi-class classifications, but this is valuable only when the tags are more specific. One particular challenge in regards to the learning process comes from the fact that if we depend solely on a semi-/supervised pipeline, we cannot know when it reaches the point where it performs better than an unsupervised one. One of our future directions involves exploring the combination of unsupervised and semi-/supervised pipelines such that they can work synchronously.

**Impact of missing values.** Through a series of experiments with the satellite experts, we found that many ML-identified anomalies are caused by missing values (see example in Fig. 8). It is known that an unsupervised ML pipeline is only able to detect unexpected temporal patterns. A missing value is one of the most frequent reasons for a particular time segment to present an unexpected "shape." Therefore, the missing value imputation strategy plays an important role in the ML pipeline. In future work, we want to seek the best way to encode missing value information as part of the time series visualization. Meanwhile, we plan to add hyperparameters to the ML pipeline that allow users to decide whether the pipeline should identify anomalies containing missing values.

## 6.3 Applicability

**Human-AI collaboration workflow usability challenges.** Although the experiments described demonstrate the success of our system, the use of ML models to facilitate decision-making may introduce some potential usability challenges [78], including lack of trust, unclear prediction targets, and difficulty reconciling human-ML disagreements. For instance, in our experiment with the satellite experts, 21% events were still marked "investigate" after 24 hours because they were hard for the experts to explain. In the long run, a growing number of such events could introduce the aforementioned ML usability challenges. Unlike more self-explanatory image or text data, time

series data can be difficult for humans to interpret. Although explainable AI (XAI) techniques for time series data [17, 27, 44, 45, 58] can alleviate potential usability challenges, they may lead to cognitive biases [67]. One promising future direction would be to formally explore which ML usability challenges exist with MTV, and investigating the best ways to integrate XAI techniques.

**Use of experiment cards.** The design of experiment cards is one of the key features we propose for maintaining team awareness during collaborative analysis. Though we did not conduct a formal study to evaluate the usefulness of experiment cards, qualitative feedback from the satellite experts confirmed the importance of this feature. In fact, experiment cards serve as an important bridge that connects every team member. P1 commented "*As a manager, this feature gives me a good overlook on the progress of each currently on-going anomaly analysis task.*". P3 confirmed "*The design of experiment cards is not only helpful for me to track the current progress but also important for me to estimate the degree of urgency.*". P5 further highlighted "*The card compose function inspires me a lot of thoughts, such as comparing different algorithms' results on the same signal.*" However, we noted that it is challenging to quickly compose a meaningful experiment card that combines as many as dozens of signals together for collaborative analysis. The current interactive table (Fig. 5-d) heavily relies on domain expertise. We consider integrating relevant signal recommendation techniques as future work.

**Real-time analysis.** Our current system is mainly used for historical data analysis. Depending on the domain need, the system can be updated every day or every week. If the system is to support real-time analysis, several questions must be addressed. How does a model know if there is a shift of data distribution in streaming data? How should the model be updated if such a shift happens? What is the optimal way to update the view with incoming data and anomalies? What changes should be made in the human-AI collaboration workflow to support real-time analysis? We aim to explore these questions as future work.

**Generalization**. The demand for time series anomaly analysis is pervasive. Although MTV was developed through close collaboration with spacecraft experts in particular, it (or parts of it) is generalizable to any application where large time series anomaly analysis is a crucial task — as are the requirements for such a system, the workflow involved, and other lessons learned from the process. Fig. 2 clarifies the scope of applications that MTV can help with by considering decision risk, time, domain expertise, technical expertise, and the role of ML. Because the current system was designed to support asynchronous collaborative analysis, it is best suited for low or medium-risk decisions that take at least a few hours. As future work, we are interested in exploring synchronous collaborative analysis to handle tasks that are of high decision risk and must be made in less time.

## 7 CONCLUSION

We have presented MTV, an interactive visual analysis system that allows multiple users to explore, investigate, and annotate multivariate time series collaboratively. The system was built to facilitate a streamlined anomaly investigation workflow, which is also summarized and presented in this work. The workflow begins with the efficient automated identification of anomalies with an end-to-end machine learning pipeline. A tailored visual interface, introduced here, allows for efficient exploration of the ML pipeline results, and features novel visualization and interaction designs that support multi-scale and multi-facet time series data exploration, as well as powerful anomaly annotation and communication. The workflow ends by closing the loop — training a semi-/supervised ML pipeline on collected annotations in order to enhance its performance. We highlight two novel algorithms — the shape-matching algorithm and the signal layout optimization algorithm — which we propose to better support the workflow. We have evaluated the speed of the shape-matching algorithm, as well as the system's ability to learn from annotations. We have conducted user studies with two groups of people: domain experts and general end-users with an

interest in analyzing time series data. Their positive feedback helps to demonstrate the effectiveness and usefulness of our system.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Subutai Ahmad, Alexander Lavin, Scott Purdy, and Zuha Agha. 2017. Unsupervised real-time anomaly detection for streaming data. *Neurocomputing* 262 (2017), 134–147.

[2] Wolfgang Aigner, Silvia Miksch, Heidrun Schumann, and Christian Tominski. 2011. *Visualization of time-oriented data.* Springer Science & Business Media.

[3] Bilal Alsallakh, Markus Bögl, Theresia Gschwandtner, Silvia Miksch, Bilal Esmael, Arghad Arnaout, Gerhard Thonhauser, and Philipp Zöllner. 2014. A visual analytics approach to segmenting and labeling multivariate time series data. In *EuroVA@ EuroVis.*

[4] Saleema Amershi and Meredith Ringel Morris. 2008. CoSearch: a system for co-located collaborative web search. In *Proceedings of the SIGCHI conference on human factors in computing systems.* 1647–1656.

[5] Jinwon An and Sungzoon Cho. 2015. Variational autoencoder based anomaly detection using reconstruction probability. *Special Lecture on IE* 2, 1 (2015), 1–18.

[6] Fabrizio Angiulli and Clara Pizzuti. 2002. Fast outlier detection in high dimensional spaces. In *European conference on principles of data mining and knowledge discovery.* Springer, 15–27.

[7] Jakob E Bardram and Steven Houben. 2018. Collaborative affordances of medical records. *Computer Supported Cooperative Work (CSCW)* 27, 1 (2018), 1–36.

[8] Jürgen Bernard, Marco Hutter, Matthias Zeppelzauer, Dieter Fellner, and Michael Sedlmair. 2017. Comparing visual-interactive labeling with active learning: An experimental study. *IEEE TVCG* 24, 1 (2017), 298–308.

[9] Jürgen Bernard, Christian Ritter, David Sessler, Matthias Zeppelzauer, Jörn Kohlhammer, and Dieter Fellner. 2017. Visual-interactive similarity search for complex objects by example of soccer player analysis. *arXiv preprint arXiv:1703.03385* (2017).

[10] Jürgen Bernard, Matthias Zeppelzauer, Michael Sedlmair, and Wolfgang Aigner. 2018. VIAL: a unified process for visual interactive labeling. *The Visual Computer* 34, 9 (2018), 1189–1207.

[11] Donald J Berndt and James Clifford. 1994. Using dynamic time warping to find patterns in time series. In *KDD workshop*, Vol. 10. Seattle, WA, USA:, 359–370.

[12] Susan E Brennan, Klaus Mueller, Greg Zelinsky, IV Ramakrishnan, David S Warren, and Arie Kaufman. 2006. Toward a multi-analyst, collaborative framework for visual analytics. In *2006 IEEE VAST.* 129–136.

[13] Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. 2000. LOF: identifying density-based local outliers. In *2000 ACM SIGMOD.* 93–104.

[14] Nan Cao, Chaoguang Lin, Qiuhan Zhu, Yu-Ru Lin, Xian Teng, and Xidao Wen. 2017. Voila: Visual anomaly detection and monitoring with streaming spatiotemporal data. *IEEE TVCG* 24, 1 (2017), 23–33.

[15] Nan Cao, Conglei Shi, Sabrina Lin, Jie Lu, Yu-Ru Lin, and Ching-Yung Lin. 2015. TargetVue: Visual analysis of anomalous user behaviors in online communication systems. *IEEE TVCG* 22, 1 (2015), 280–289.

[16] Varun Chandola, Arindam Banerjee, and Vipin Kumar. 2009. Anomaly detection: A survey. *ACM computing surveys (CSUR)* 41, 3 (2009), 15.

[17] Furui Cheng, Dongyu Liu, Fan Du, Yanna Lin, Alexandra Zytek, Haomin Li, Huamin Qu, and Kalyan Veeramachaneni. 2021. VBridge: Connecting the Dots Between Features and Data to Explain Healthcare Models. *IEEE Transactions on Visualization and Computer Graphics* (2021).

[18] Tshering Dema, Margot Brereton, Jessica L. Cappadonna, Paul Roe, Anthony Truskinger, and Jinglan Zhang. 2017. Collaborative Exploration and Sensemaking of Big Environmental Sound Data. *Computer Supported Cooperative Work (CSCW)* 26, 4-6 (Dec. 2017), 693–731.

[19] Johannes Fuchs, Fabian Fischer, Florian Mansmann, Enrico Bertini, and Petra Isenberg. 2013. Evaluation of alternative glyph designs for time series data in a small multiple setting. In *Proceedings of the SIGCHI conference on human factors in computing systems.* 3237–3246.

[20] Vagner Luiz Gava, Mauro de Mesquita Spinola, Antonio Carlos Tonini, and José Cardenas Medina. 2012. The 3c cooperation model applied to the classical requirement analysis. *JISTEM-Journal of Information Systems and Technology Management* 9, 2 (2012), 235–264.

[21] Alexander Geiger, Dongyu Liu, Sarah Alnegheimish, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. 2020. TadGAN: Time Series Anomaly Detection Using Generative Adversarial Networks. *arXiv preprint arXiv:2009.07769* (2020).

[22] Anna Gogolou, Theophanis Tsandilas, Themis Palpanas, and Anastasia Bezerianos. 2018. Comparing similarity perception in time series visualizations. *IEEE TVCG* 25, 1 (2018), 523–533.

[23] Markus Goldstein and Seiichi Uchida. 2016. A comparative evaluation of unsupervised anomaly detection algorithms for multivariate data. *PloS one* 11, 4 (2016), e0152173.

[24] Nitesh Goyal and Susan R Fussell. 2016. Effects of sensemaking translucence on distributed collaborative analysis. In *Proceedings of the 19th ACM Conference on Computer-Supported Cooperative Work & Social Computing*. 288–302.

[25] Riyaz Ahamed Ariyaluran Habeeb, Fariza Nasaruddin, Abdullah Gani, Ibrahim Abaker Targio Hashem, Ejaz Ahmed, and Muhammad Imran. 2019. Real-time big data processing for anomaly detection: A survey. *International Journal of Information Management* 45 (2019), 289–307.

[26] Amir Hossein Hajizadeh, Melanie Tory, and Rock Leung. 2013. Supporting awareness through collaborative brushing and linking of tabular data. *IEEE TVCG* 19, 12 (2013), 2189–2197.

[27] Jonathan J Harris, Ching-Hua Chen, and Mohammed J Zaki. 2020. A Framework for Generating Explanations from Temporal Personal Health Data. *arXiv preprint arXiv:2003.09530* (2020).

[28] Jeffrey Heer and Maneesh Agrawala. 2008. Design considerations for collaborative visual analytics. *Information visualization* 7, 1 (2008), 49–62.

[29] Jeffrey Heer, Frank Van Ham, Sheelagh Carpendale, Chris Weaver, and Petra Isenberg. 2008. Creation and collaboration: Engaging new audiences for information visualization. In *Information visualization*. Springer, 92–133.

[30] Jeffrey Heer, Fernanda B Viégas, and Martin Wattenberg. 2007. Voyagers and voyeurs: supporting asynchronous collaborative information visualization. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. 1029–1038.

[31] Florian Heimerl, Steffen Koch, Harald Bosch, and Thomas Ertl. 2012. Visual classifier training for text document retrieval. *IEEE TVCG* 18, 12 (2012), 2839–2848.

[32] Victoria Hodge and Jim Austin. 2004. A survey of outlier detection methodologies. *Artificial intelligence review* 22, 2 (2004), 85–126.

[33] Benjamin Höferlin, Rudolf Netzel, Markus Höferlin, Daniel Weiskopf, and Gunther Heidemann. 2012. Inter-active learning of ad-hoc classifiers for video visual analytics. In *2012 IEEE VAST*. 23–32.

[34] Sungsoo Hong, Minhyang Suh, Nathalie Henry Riche, Jooyoung Lee, Juho Kim, and Mark Zachry. 2018. Collaborative dynamic queries: Supporting distributed small group decision-making. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. 1–12.

[35] Sungsoo Hong, Minhyang Suh, Tae Soo Kim, Irina Smoke, Sangwha Sien, Janet Ng, Mark Zachry, and Juho Kim. 2019. Design for Collaborative Information-Seeking: Understanding User Challenges and Deploying Collaborative Dynamic Queries. *Proceedings of the ACM on Human-Computer Interaction* 3, CSCW (2019), 1–24.

[36] Kyle Hundman, Valentino Constantinou, Christopher Laporte, Ian Colwell, and Tom Soderstrom. 2018. Detecting spacecraft anomalies using lstms and nonparametric dynamic thresholding. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 387–395.

[37] Petra Isenberg and Sheelagh Carpendale. 2007. Interactive tree comparison for co-located collaborative information visualization. *IEEE Transactions on Visualization and Computer Graphics* 13, 6 (2007), 1232–1239.

[38] Petra Isenberg, Niklas Elmqvist, Jean Scholtz, Daniel Cernea, Kwan-Liu Ma, and Hans Hagen. 2011. Collaborative visualization: Definition, challenges, and research agenda. *Information Visualization* 10, 4 (2011), 310–326.

[39] Petra Isenberg and Danyel Fisher. 2009. Collaborative brushing and linking for co-located visual analytics of document collections. In *Computer Graphics Forum*, Vol. 28. Wiley Online Library, 1031–1038.

[40] Petra Isenberg, Danyel Fisher, Sharoda A Paul, Meredith Ringel Morris, Kori Inkpen, and Mary Czerwinski. 2011. Co-located collaborative visual analytics around a tabletop display. *IEEE TVCG* 18, 5 (2011), 689–702.

[41] F. Itakura. 1975. Minimum prediction residual principle applied to speech recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing* 23, 1 (1975), 67–72.

[42] Waqas Javed, Bryan McDonnel, and Niklas Elmqvist. 2010. Graphical perception of multiple time series. *IEEE TVCG* 16, 6 (2010), 927–934.

[43] Robert Kincaid and Heidi Lam. 2006. Line graph explorer: scalable display of line graphs using focus+ context. In *Proceedings of the working conference on Advanced visual interfaces*. 404–411.

[44] Jiwei Li, Will Monroe, and Dan Jurafsky. 2016. Understanding neural networks through representation erasure. *arXiv preprint arXiv:1612.08220* (2016).

[45] Dongyu Liu, Weiwei Cui, Kai Jin, Yuxiao Guo, and Huamin Qu. 2018. Deeptracker: Visualizing the training process of convolutional neural networks. *ACM Transactions on Intelligent Systems and Technology (TIST)* 10, 1 (2018), 1–25.

[46] Weichen Liu, Sijia Xiao, Jacob T. Browne, Ming Yang, and Steven P. Dow. 2018. ConsensUs: Supporting Multi-Criteria Group Decisions by Visualizing Points of Disagreement. 1, 1 (2018).   https://doi.org/10.1145/3159649

[47] Thomas Ludwig, Tino Hilbert, and Volkmar Pipek. 2015. Collaborative visualization for supporting the analysis of mobile device data. In *2015 ECSCW*. Springer, 305–316.

[48] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, Nov (2008), 2579–2605.

[49] Narges Mahyar and Melanie Tory. 2014. Supporting communication and coordination in collaborative sensemaking. *IEEE TVCG* 20, 12 (2014), 1633–1642.

[50] Peter McLachlan, Tamara Munzner, Eleftherios Koutsofios, and Stephen North. 2008. LiveRAC: interactive visual exploration of system management time-series data. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 1483–1492.

[51] Milagros Miceli, Martin Schuessler, and Tianling Yang. 2020. Between Subjectivity and Imposition: Power Dynamics in Data Annotation for Computer Vision. *Computer Supported Cooperative Work (CSCW)* 4, 2 (2020), 1–25.

[52] Meredith Ringel Morris. 2013. Collaborative search revisited. In *Proceedings of the 2013 conference on Computer supported cooperative work*. 1181–1192.

[53] Meredith Ringel Morris and Eric Horvitz. 2007. SearchTogether: an interface for collaborative web search. In *Proceedings of the 20th annual ACM symposium on User interface software and technology*. 3–12.

[54] Eduardo HM Pena, Marcos VO de Assis, and Mario Lemes Proença. 2013. Anomaly detection using forecasting methods arima and hwds. In *2013 IEEE International Conference of the Chilean Computer Science Society (SCCC)*. 63–66.

[55] Huamin Qu, Wing-Yi Chan, Anbang Xu, Kai-Lun Chung, Kai-Hon Lau, and Ping Guo. 2007. Visual analysis of the air pollution problem in Hong Kong. *IEEE TVCG* 13, 6 (2007), 1408–1415.

[56] H. Sakoe and S. Chiba. 1978. Dynamic programming algorithm optimization for spoken word recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing* 26, 1 (1978), 43–49.

[57] Advait Sarkar, Martin Spott, Alan F Blackwell, and Mateja Jamnik. [n.d.]. Visual discovery and model-driven explanation of time series patterns. In *2016 IEEE Symposium on Visual Languages and Human-Centric Computing*. 78–86.

[58] Udo Schlegel, Hiba Arnout, Mennatallah El-Assady, Daniela Oelke, and Daniel A Keim. 2019. Towards a rigorous evaluation of XAI Methods on Time Series. *arXiv preprint arXiv:1909.07082* (2019).

[59] Burr Settles. 2011. Closing the loop: Fast, interactive semi-supervised annotation with queries on features and instances. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*. 1467–1478.

[60] Chirag Shah. 2012. *Collaborative information seeking: The art and science of making the whole greater than the sum of all*. Vol. 34. Springer Science & Business Media.

[61] Ben Shneiderman. 2003. The eyes have it: A task by data type taxonomy for information visualizations. In *The craft of information visualization*. Elsevier, 364–371.

[62] Micah J Smith, Carles Sala, James Max Kanter, and Kalyan Veeramachaneni. 2020. The machine learning bazaar: Harnessing the ML ecosystem for effective system development. In *2020 ACM SIGMOD*. 785–800.

[63] Megan K Torkildson, Kate Starbird, and Cecilia Aragon. 2014. Analysis and visualization of sentiment and emotion on crisis tweets. In *International conference on cooperative design, visualization and engineering*. Springer, 64–67.

[64] Jarke J Van Wijk and Edward R Van Selow. 1999. Cluster and calendar based visualization of time series data. In *1999 InfoVis*. 4–9.

[65] Fernanda B Viegas and Martin Wattenberg. 2006. Communication-minded visualization: A call to action. *IBM Systems Journal* 45, 4 (2006), 801.

[66] MK Vijaymeena and K Kavitha. 2016. A survey on similarity measures in text mining. *Machine Learning and Applications: An International Journal* 3, 2 (2016), 19–28.

[67] Danding Wang, Qian Yang, Ashraf Abdul, and Brian Y Lim. 2019. Designing theory-driven user-centric explainable AI. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 1–15.

[68] Colin Ware. 2019. *Information visualization: perception for design*. Morgan Kaufmann.

[69] Marc Weber, Marc Alexa, and Wolfgang Müller. 2001. Visualizing time-series on spirals. In *2001 InfoVis*, Vol. 1. 7–14.

[70] Wesley Willett, Jeffrey Heer, Joseph Hellerstein, and Maneesh Agrawala. 2011. CommentSpace: structured support for collaborative visual analysis. In *Proceedings of the SIGCHI conference on Human Factors in Computing Systems*. 3131–3140.

[71] Jing Xia, Jieqiong Zhao, Isaac Sheeley, Joseph Christopher, Qiaoying Wang, Chen Guo, Jiawei Zhang, David S Ebert, Yingjie Victor Chen, and Zhenyu Cheryl Qian. 2014. AnnotatedTimeTree: Visualization and annotation of news text and other heterogeneous document collections. In *2014 IEEE VAST*. 337–338.

[72] Cong Xie, Wei Xu, and Klaus Mueller. 2018. A visual analytics framework for the detection of anomalous call stack trees in high performance computing applications. *IEEE TVCG* 25, 1 (2018), 215–224.

[73]  Ke Xu, Shunan Guo, Nan Cao, David Gotz, Aiwen Xu, Huamin Qu, Zhenjie Yao, and Yixin Chen. 2018. Ecglens: Interactive visual exploration of large scale ecg data for arrhythmia detection. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 1–12.

[74]  Ke Xu, Yun Wang, Leni Yang, Yifang Wang, Bo Qiao, Si Qin, Yong Xu, Haidong Zhang, and Huamin Qu. 2019. CloudDet: Interactive Visual Analysis of Anomalous Performances in Cloud Computing Systems. *IEEE TVCG* 26, 1 (2019), 1107–1117.

[75]  Jian Zhao, Michael Glueck, Petra Isenberg, Fanny Chevalier, and Azam Khan. 2017. Supporting handoff in asynchronous collaborative sensemaking using knowledge-transfer graphs. *IEEE TVCG* 24, 1 (2017), 340–350.

[76]  Dequan Zheng, Fenghuan Li, and Tiejun Zhao. 2016. Self-adaptive statistical process control for anomaly detection in time series. *Expert Systems with Applications* 57 (2016), 324–336.

[77]  Bin Zhou, Shenghua Liu, Bryan Hooi, Xueqi Cheng, and Jing Ye. 2019. BeatGAN: Anomalous Rhythm Detection using Adversarially Generated Time Series. In *IJCAI*. 4433–4439.

[78]  Alexandra Zytek, Dongyu Liu, Rhema Vaithianathan, and Kalyan Veeramachaneni. 2021. Sibyl: Understanding and Addressing the Usability Challenges of Machine Learning In High-Stakes Decision Making. *IEEE Transactions on Visualization and Computer Graphics* (2021).