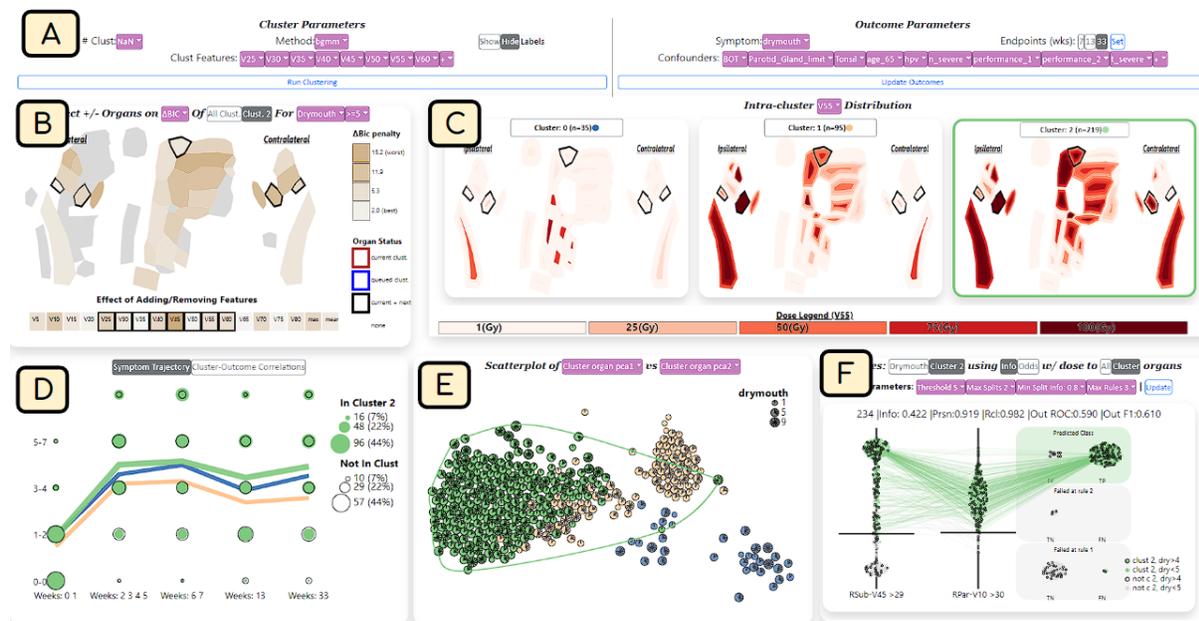


# DASS Good: Explainable Data Mining of Spatial Cohort Data

A. Wentzel<sup>1</sup>, C. Floricel<sup>1</sup>, G. Canahuate<sup>2</sup>, M.A. Naser<sup>3</sup>, A.S. Mohamed<sup>3</sup>, C.D. Fuller<sup>3</sup>, L. van Dijk<sup>3</sup> and G.E. Marai<sup>1</sup>

<sup>1</sup> University of Illinois Chicago, Electronic Visualization Lab <sup>2</sup>University of Iowa <sup>3</sup>University of Texas MD Anderson Cancer Center



**Figure 1:** DASS interactive model building for head and neck cancer. A) Control panel for changing cluster parameters and the desired outcome. B) Additive Effect panel showing the effect of changing cluster features. C) Intra-cluster dose distribution plot. D) Outcome plot showing the symptom ratings of patients over time within each cluster. E) Stylized scatterplot showing cohort projections. F) Rule builder view, showing a rule-based mimic model that predicts patients in the selected cluster.

## Abstract

Developing applicable clinical machine learning models is a difficult task when the data includes spatial information, for example, radiation dose distributions across adjacent organs at risk. We describe the co-design of a modeling system, DASS, to support the hybrid human-machine development and validation of predictive models for estimating long-term toxicities related to radiotherapy doses in head and neck cancer patients. Developed in collaboration with domain experts in oncology and data mining, DASS incorporates human-in-the-loop visual steering, spatial data, and explainable AI to augment domain knowledge with automatic data mining. We demonstrate DASS with the development of two practical clinical stratification models and report feedback from domain experts. Finally, we describe the design lessons learned from this collaborative experience.

## 1 Introduction

Precision radiotherapy (RT) is a medical paradigm that seeks to personalize cancer RT and care for an individual patient, based on data from cohorts of similar patients. Because for many site-specific cancers, the treatment depends on the location and spread of the disease, modern approaches to precision RT aim to leverage spatial patient-specific information such as anatomical data drawn from CT scans [WHL\*19, WHvD\*20]. In conjunction with the co-

hort data, this information can then be used to improve patient outcomes such as survival or quality of life after treatment.

In this context, machine learning (ML) models are powerful tools for stratifying the cohort data in meaningful ways, for example into patient groups at high-risk versus low-risk of developing treatment-related symptoms. However, developing applicable clinical ML models for patient stratification is difficult when the data includes spatial information, for example, radiation dose distributions across adjacent organs at risk. In addition, while ML approaches of-

ten work well with large oncology data, automated model-building approaches using smaller cohorts often perform poorly when deployed in practice [MK22]. Furthermore, prediction using treatment plans and qualitative outcomes such as symptom ratings is particularly difficult. This results in simpler models that may underperform or complex models that are very likely to overfit. With advancements in explainable AI techniques, we can better probe models and iteratively find ways of improving models that properly leverage domain knowledge, helping us avoid issues with poor generalization and overfitting, while improving on standard statistical approaches. These combined issues make RT cohort modeling well-suited for a human-machine mixed-initiative system.

In this work, we present a visual steering approach for creating patient stratifications of head and neck cancer (HNC) patients based on 3-dimensional dose distributions to organs-at-risk, to separate patients at high risk of experiencing long-term side effects. Unlike the current state of the art, our approach supports interactively exploring and visualizing high-dimensional spatial dose distributions, the temporal analysis of RT cohort data, access to both individual patient data and patient distribution within a cluster, constructing unsupervised rule models to help explain the clusters, and iteratively refining and exploring parameters to create actionable stratifications. We implement this approach in *Dose Analytics and Symptom Stratifier* (DASS), a visual computing system designed to allow for the development and exploration of patient stratifications according to different symptoms of interest. We describe two case studies of applying DASS and show how it has been used to improve existing outcome models. Finally, we provide design lessons gained through this collaborative visual steering design.

## 2 Background

Head and neck oncology has seen large increases in patient survival due to a shift from smoking-driven tumors to less aggressive HPV-driven tumors. This increase in survival has resulted in a shift in priorities towards increasing the quality of life of patients: radiation to organs near the primary tumor during treatment can lead to tissue damage, resulting in long-term side effects [LDVdL\*08, ELB\*91]. Predicting when symptoms driven by spatial tissue damage occur is thus an understudied area of interest to oncologists, as it can help identify better treatment guidelines.

When performing the initial diagnosis, oncologists rely on patient history and clinical staging that rank the size and spread of the tumor [OHS\*16] to determine the method of treatment to optimize patient survival. However, after the treatment methodology is established, predictive models are needed to identify patients that may need preventative treatment for serious side effects. A diagram of the clinical workflow is available in the supplementary material (Figure A1).

In particular, predicting tissue damage from radiation therapy in head and neck cancer (HNC) patients is a challenge due to the high number of treatment parameters and high number of organs that may factor into side effects. For example, drymouth is often caused by radiation damage to a subset of the salivary glands. Identifying when failure may occur is a difficult modeling task, in which one needs to consider the glands as a spatially interrelated system, as some may compensate for damage to other glands. Additionally, each organ may have a separate non-linear response to the radiation

dose over time, and symptom severity varies throughout treatment. Furthermore, the large numbers of HNC patients in a cohort and the dimensionality of the data pose a challenge in terms of visual analysis. Finally, human modelers also require access to individual patient data, as well as to the patient distribution within a cluster to make informed inferences about patient outcomes.

## 3 Related Work

### 3.1 Visual Analysis of Cohort Data

Several applications of visual analysis have focused on different algorithmic approaches for clustering patients [MKKW12, MAM20, WHvD\*20]. Visualization tools often extend these approaches by allowing human-in-the-loop analysis to identify sub-cohorts [ZGP15, KPS16, BSM\*15]. Other systems have focused on comparison of cohorts to discover differences in disease progression [MDM\*15], genetics [GGC\*17], cancer treatment disparities [STA\*22], but, unlike our work, these systems do not focus on model building.

Many systems use clustering [MV15] and dimensionality reduction [ODH\*07, ENBD08] on key features to guide explorations over high-dimensionality data. Some tools have looked at visual analytics for creating clusters with unstructured health data [KEV\*17, CD18, GNDV\*17], while other systems incorporate temporal clustering methods [ZMP\*21, ZMW\*20a, FNB\*21, WMH\*21, GXZ\*17]. However, these systems do not attempt to incorporate spatial information in their clustering models, as we do. Additionally, none of these systems link detailed treatment plans to qualitative patient outcomes in the cohorts, as we do.

### 3.2 Visualization of Medical Image Data

Work in visual computing with medical imaging often focuses on linking spatial features to external variables to support exploration for domain experts. Early work focused on visualizing spatial imaging data with open source tools (MITK [VWW\*04]) and introduced integration of spatial and non-spatial linked views [GRW\*00].

Specialized approaches have been developed to explore cohort features in other domains such as tissue imaging [FYTL18, JKW\*22, WKN\*12], neuroscience [JBB\*08, JBF\*20, ASO\*16, MPL\*18], and lumbar spine features [CLL\*21, KOJL\*14].

Focusing on cohorts of RT data, BladderRunner [RCMA\*18] visualized cohorts of prostate cancer patients which used a mixture of T-SNE and Gaussian mean-shift clustering to group patients based on bladder shape. VAPOR [FGM\*20] extended their work to consider RT-induced treatment toxicity. Other work has extended these results to explore uncertainty in RT data for visual analysis [GCMM\*19, RPHL14] and predictive models [FMCM\*21]. However, these approaches do not deal with HNC oncology treatment, which has more complex treatment and symptom patterns but lower temporal variability.

Previous HNC work has used spatial data to cluster patients based on tumor spread to lymph nodes [LWE\*20]. Many techniques rely on simplified representations of anatomical data to allow for better analysis of high-dimensional data [WCVD\*20, WHL\*19, KOJL\*14, RCMA\*18]. While these works often deal with feature engineering, none of them focus on directly altering

the model in parallel with the visual analysis, as we do. Additionally, we uniquely provide tools for validating the feasibility of the underlying model’s logic and embedding anatomical data directly into the system.

### 3.3 Visual Steering and Interactive Machine Learning

In the medical domain, several projects have developed visualization systems around the workflows of clinical model builders and biostatisticians with a focus on regression models [DvVH\*19]. Raidou et. al [RvdHD\*15, RCMM\*16] proposed a tool for visual analysis of regression-based Tissue Complication Probability models, with a focus on uncertainty. However, these approaches do not focus on clustering or stratification models, as we do.

Other work has focused on actionable explanations for pre-built models for clinicians, such as normal tissue complication models [ZHT\*13], binary classifiers [CMQ20], case-based reasoning [MMB\*18, MXC\*19], and black box models [CLD\*21]. For explainable AI, DrugExplorer [WHC\*22] proposed a model for user-centered XAI alongside a system for exploring graph-neural-networks for drug repurposing. However, none of these approaches tackle iterative probing and model development, or capturing spatial information in their data, as we do.

Additionally, our work uses interactive rule mining to help explain the clusters. Many systems have worked on aggregated visualization of rules [SMS\*22, TM03b, TM03a, vdEvW11, MP13, XSFM11, YNB21], and used interactive rule mining to approximate more complex models [MQB18]. Our approach differs from these in that we include a novel rule mining algorithm focused on matching clinical use cases, along with a novel visual encoding that allows for interactive parameter tuning.

## 4 Methods

The DASS design is rooted in our earlier experience with clinical stratification models that relied on forward search for feature selection for clustering [WHvD\*20]. Fully automated parameter searches yielded models that performed well on a single performance metric. However, when the clusters were inspected by clinical collaborators, they would often find issues with the organs used, such as organs that are completely unrelated to the outcome, or smaller organs that they felt should be included. Thus, we introduced a human-in-the-loop forward search directly into our front-end alongside model explanations to help improve the process of iterating on our clusters.

User-guided search has two additional benefits. First, our clinician collaborators wished to specify desired characteristics of the models, which led to a need to explore multiple alternative outcomes or starting points based on these desired characteristics. Second, collaborator input is required when balancing model performance, the feasibility of the organs considered, and the number of organs considered. For example, we found that in one model, including both the soft and hard palate had identical effects on the outcome. Thus, the decision came down to the clinicians, who helped us identify which one was of more clinical importance.

Furthermore, in previous work, we attempted to find clusters through hyperparameter search or using predefined cluster features. However, we found that neither approach performed well.

Automatic feature selection led to clusters that focus on organs that served as positional *indicator features*, such as the oral cavity [WHvD\*20], but are not causally linked to outcomes and resulted in model explanations that are not well-received. Notably, we found that the brainstem and brachial plexus nerves often appeared as predictors, despite clinicians noting that neither can be associated with any of the outcomes being predicted. Such models work well, but lack causal plausibility which hinders adoption and cannot be generalized to treatment guidelines. The DASS design specifically addresses these problems through its back-end and front-end.

### 4.1 Data

Data were collected from a cohort of 349 HNC patients treated at the MD Anderson Cancer Center using Radiation Therapy, with or without chemotherapy, using a 7-week treatment course. We consider three types of data: spatial dosimetric data taken from the patient’s treatment plan; unstructured clinical data taken from the patient’s health record; and temporal information on the patient’s self-reported side effects taken during and after treatment. All values are positive ordinal values. Symptom ratings for individuals are discrete, while dose values are continuous.

Diagnostic images were taken at the time of diagnosis, and 40 organs of interest were segmented from these images and considered in the treatment plan. Dose treatment plans were extracted for each organ of each patient. We include 3-dimensional information on the cumulative dose received by each organ during treatment. We use the notation “VX” to denote the maximum dose that penetrates X% of the organ. For each organ, we consider the V5-V95 range in increments of 5, as well as the mean and maximum dose.

For outcomes, patients were asked to fill out an MD Anderson Symptom Inventory (MDASI) questionnaire [RMC\*07]. This inventory includes self-reported symptoms for 28 different items, such as drymouth and pain on a scale of 1-10. We also include secondary variables that may be used as confounders in the patient outcomes taken from electronic health record data, which we generally treat as binary confounding variables.

### 4.2 Collaboration

Our work was done as part of an ongoing collaboration between data scientists and research oncologists at three US sites. DASS was commissioned to serve first and foremost the needs of the model builders, but to also facilitate clinician input and feedback on the models. Remote meetings were held weekly, during which we would get feedback on designs, and update project goals based on feedback and current results. Examples of prototypes during this phase are included in the supplement Appendix B.

We followed an Activity-Centered Design (ACD) process [Mar17], which is a methodology conceived to better support designing for domain experts by focusing on existing user workflows and activities. The approach has higher success rates in interdisciplinary settings than Human-Centered Design (63% versus 25%) [Mar17]. We focused on the workflows around the development of clinically applicable models, as well as the associated data analysis and verification required to validate and publish the results.

### 4.3 Task Analysis

**Modeling Requirements** The goal of our project was to aid in the development of an interpretable decision-support tool for clinicians to help identify HNC patients at high risk of long-term severe (self-reported rating  $> 4$  on a 10-point scale) side effects due to radiation damage. We focus on HNC patients as the sensitivity of organs in the head and neck makes detection of quality of life measures in survivors a difficult, under-explored application. Our collaborators were specifically interested in a model that could improve on existing clinical systems by incorporating sets of related organs that together support specific functions, and thus should be treated as a system.

Our system was designed to be used for asymmetric collaborative analysis, which would be handled by model-builders with expertise in the underlying algorithms, with clinicians providing input and feedback. Therefore, we identified requirements for the models themselves, as well as the steps needed to create and validate each model. For our models, we derived the following requirements:

*Actionable:* Usable in a practical setting. In a typical workflow, clinicians use risk stratifications that rank a patient's risk of survival, which are then integrated into a holistic treatment plan. As such, we require that our models output a simple ranking for each patient, as well as insights that are usable without access to the models. Access to individual patient data, as well as the patient distribution in each cluster, in terms of both doses and symptoms, was necessary.

*Plausible:* Generalize well to a real-world setting. The underlying features that lead to a patient being classified as high-risk must be easy to understand in their spatial context. The models must also place patients in the high-risk group because they received a high dose to a specific set of organs, and the set of organs considered must be mechanistically linked to the the outcome of interest.

*Transparent:* Be easily probed, assess the plausibility of the models, and identify edge cases in the models. We also needed to be able to demonstrate the plausibility of the models and explain its internal logic to readers with a clinical background.

Based on these requirements, we designed a dose-based stratification methodology that clustered 2D dose distributions to a set of organs and used the resulting patient clusters as a proxy for patient risk. Our visual front-end is designed around visual steering, which uses information scent and visual cues to guide our team through the process of selecting, validating, and refining the range of potential parameters for the models to balance different performance metrics and model plausibility. Because this task requires significant knowledge of the models when adjusting parameters, our interactive system is designed to be used directly by models builders and visual computer experts, with encodings designed to allow model builders to communicate intermediate results to clinicians and domain experts.

Through a series of iterative sessions where we developed models and discussed them with our collaborators, we identified the following Activities and Tasks for our visual interface:

- **A1** - Given a symptom, find optimal cluster parameters
  - **T1** Find organs causally related to the symptom of interest.

- **T2** Identify a window in the dose-volume histogram that best stratifies the cohort.
- **T3** Validate a choice of clustering algorithm and parameters
- **A2** - Validate that the logic of a model is causal and plausible
  - **T4** Examine the dose distribution of each cluster and where the doses differ.
  - **T5** Verify if the cluster with the highest symptom risk also has higher doses to the organs used in the clustering.
  - **T6** Identify confounders that may impact risk prediction.
  - **T7** Validate the predictive accuracy of the clusters.
- **A3** - Examine and explain individual clusters
  - **T8** Identify the organ doses that most distinguish each cluster.
  - **T9** Evaluate differences in symptom trajectory between clusters over time.

A1 deals with the development of a models, while activities A2 and A3 help to quantify the models and provide feedback to improve the parameters in A1. A2 is a requirement for clinical publishable findings, while A3 is important for identifying any insights that can be drawn from the final model. For example, once a model is validated, finding that the high-risk cluster for taste dysfunction tends to have a very high maximum dose to the tongue may indicate that future work should investigate the effect of tongue dose on outcomes in more detail.

### 4.4 Back-end Algorithms

**Modeling.** DASS allows selecting from a range of clustering algorithms: K-nearest-neighbors, Hierarchical clustering, spectral clustering, and a Gaussian Mixture Model. After several iterations, we converged to a Bayesian variant of a Gaussian mixture model for all cluster outcomes. Once a set of organs and a dose-volume histogram (DVH) is identified, these features are encoded as a vector for each patient of size  $\#organs * window-size$ . Patient vectors are clustered, which are ranked based on the sum of the mean doses to each organ included in the cluster. Ideally, this will result in the highest rank cluster (high dose) being the most correlated with the outcome.

To evaluate the resulting models, we also need to specify a symptom and time point to use as the outcome of interest. We then convert ratings to a binary outcome using a severity threshold. After discussion with our collaborators, the default was a symptom rating above 4 out of 10 at 6 months after treatment.

Once our clusters and outcomes are identified, we perform multivariate correlation analysis using a likelihood ratio test (LRT) to assess the correlation between each cluster and the outcome of interest, using a set of clinical confounders interactively specified.

From this, we can calculate an odds-ratio and statistical significance p-value for each cluster, as well as the Bayesian Information Change (BIC) [KK08]. BIC and AIC are estimates of the goodness of fit of a model that include a penalty for the number of variables considered, in order to prevent overfitting, where lower scores indicate better fits [KK08]. For BIC, reductions in score relative to a baseline model of at least 2 indicate reasonable evidence, while reductions of at least 6 indicate "strong" evidence of improvement [Raf95]. This provides a set of different metrics for assessing the cluster quality in terms of stratifying the cohort.

In addition, to assessing the quality of the current clustering, we provide a forward search in which we alter the existing cluster parameters by adding or removing either a single organ or a single feature from the dose-volume histogram window. We then re-cluster the cohort, and evaluate the new p-value, AIC, and BIC for the new clusters, relative to that of the existing cluster. These metrics are used to provide information scent for users when performing a forward search of the data.

**Rule Mining.** To help explain the clusters, we designed a constrained rule mining algorithm and used it to generate a set of dose thresholds that work as a classifier. Our algorithm looks for splits among all dose features in the dataset to find a set that maximizes the mutual information between the splits and a binary outcome. This algorithm is designed to approximate standard rule mining, with the following additional constraints so that the results approximate the rules used by clinicians when specifying dose thresholds: 1) Monotonicity – the high-probability subset for each split in the data must either always be the group above or always in the group below the threshold; 2) Minimalism – The algorithm can only use one dose-feature for each organ; 3) Informative – each “rule” in the ruleset must have a minimal predictive value (user-set) on its own.

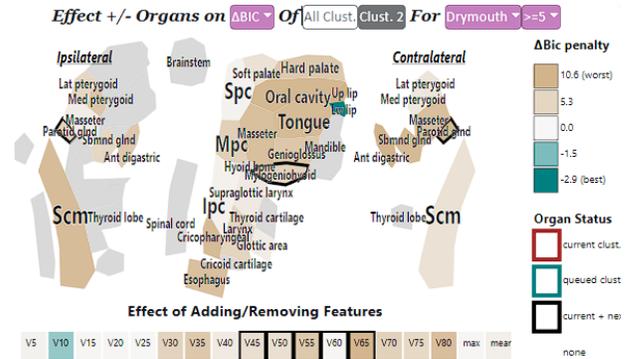
Specifically, the algorithm works as follows: 1) we calculate the mutual information gain between each feature split within each ROI (e.g. V40 to the Tongue > 40) and the binary outcome of interest; 2) of the resulting splits, we select the k most important splits; 3) for each of the k best rules, we test combinations of all other splits in step 1 that do not share the same ROI, and calculate the new mutual information gain of the combined rules. Rules are combined using the AND operator (i.e. the patients must satisfy all rules); 4) steps 2-3 are repeated until no improvement is seen in the mutual information gain. To speed up the algorithm, pruning parameters used to speed up the search can be adjusted in the interface.

**Implementation** All data pre-processing and modeling is done using Python with NumPy, Pandas, and Flask for the back-end. Clustering and dimensionality reduction is performed using the scikit-learn package while statistical tests use the statsmodels package. Our system frontend is implemented using React and D3.js.

## 5 Front-end Design

The DASS front-end (Fig. 1) is composed of 6 panels: a cluster dose view (Fig. 1-B) that shows the within-organ dose distribution for each cluster (A2), an additive effects view (Fig. 1-C) that shows the estimated impact of adding or removing features from the cluster on the specified outcome (A1), an outcome view (Fig. 1-D) that shows the different symptom ratings over time for each cluster (A2-A3), a configurable scatterplot view (Sec. 1-E) that shows a 2D projection of all the patients in either the dose or outcome space (A2), a rule view (Fig. 1-F) that shows a set of dose thresholds that best separates a cluster of interest (A1-A3), and a control panel (Fig. 1-A) that allows users to specify the cluster parameters and outcomes of interest. We arrived at this design following a parallel prototyping process, with multiple design alternatives and repeated feedback. This process is illustrated in the supplemental materials.

To better support the analytical workflow, we use a categorical color scale for cluster membership. Analysts can select a specific cluster, which is used to populate the temporal outcome and rule



**Figure 2:** Additive Effects encoding showing a heat map of the organs and dose-features used in clustering. Color encodes the goodness of fit effect of adding (no or teal outline) or removing (dark black or brown outline) features to the clustering.

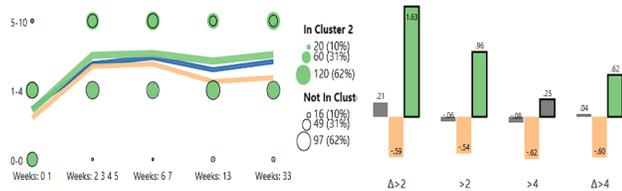
views, and brush in all other linked views. By default, DASS automatically selects for brushing the highest dose cluster, as this cluster was typically of the most interest to our clinicians.

### 5.1 Visual Scaffolding

When dealing with organ data, understanding the relative position of each organ is essential for analysis of the relationships between organs and side-effects. Specifically, dose values are correlated with location, and it is important to identify situations where organs may be linked to toxicities due to their centrality and proximity to nearby organs rather than being directly causally linked.

In previous work, we represented the set of organs as a stylized plot showing each organ as a plot in 3 dimensions [WHL\*19]. However, we felt that this representation was limited in its usefulness, as it is difficult to identify organs that may be smaller and clustered together, but may be functionally important, such as salivary glands and smaller organs in the neck. Previous work has also shown that 2-dimensional maps of anatomical regions work well, and work well with clinicians who are typically trained to work with image slices and 2-dimensional anatomical drawings [WCVD\*20]. Expanding on this, we created a 2-dimensional representation of 45 organs used in our dataset based on existing anatomical drawings [FH15].

We then divided up the organs in the head into unilateral organs that sit along the mid-sagittal plane (e.g. tongue), and those that exist as a pair of organs on each side of the mid-sagittal plane (e.g. eyes), which are further subdivided into those on the same side as the primary tumor (ipsilateral side) and those on the opposite side of the primary tumor (contralateral side). This gives us three “groups” of organs along the center axis. For each region, we took tracings around organs of interest using multiple anatomical cross-sections. We then overlaid all drawings, added in missing regions such as the spinal cord, and manually adjusted each contour to avoid overlap and regularize the size of each region. Adjustments were also made to ensure that regions were reasonably concave so that color gradients were visible. A diagram of the final drawing with all regions labeled is available in the supplementary materials.



**Figure 3:** (Left) Plots showing the symptom ratings over time from the start of treatment for the specified symptom of interest, broken up by cluster. Circular markers encode the percentage of patients that experience a symptom at each level and time point, and help us estimate a patient’s relative risk. Line charts show average ratings for each symptom. (Right) Bar chart showing the results of multi-variate correlation tests for the clusters at different thresholds.

## 5.2 Additive Effects Panel

When working on model development (A1) our main task is to identify a set of organs to cluster once our desired outcome has been specified (Fig. 2). In this panel, we provide a forward search to estimate the effect of adding (for features not in the current clusters) or removing (for features in the current clusters) different organs or features from the clustering space on model performance (Sec. 4.4). We chose a beige-white-teal color scheme as we wanted to de-emphasize uninteresting (negative) results while still capturing the divergent nature of the results. Thus, we used beige as it has lower perceptual salience than the rest of DASS.

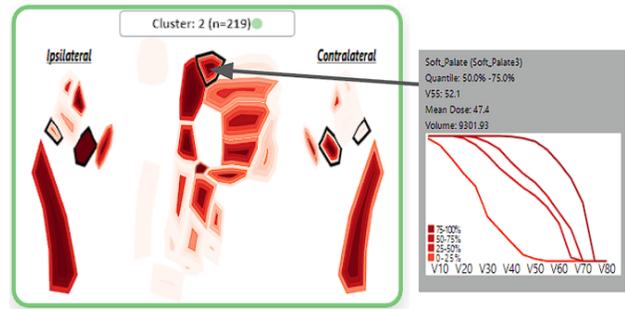
Since model developers may be interested in balancing performance between multiple outcomes, we allow choosing which information metric is used to encode color: BIC, AIC, or the t-statistic—which we report as a change in p-value, as well as the inputs to the LRT test, and the threshold used to rank an outcome as “severe”.

Alternative designs relied on variations of heat-map and bar charts with effect sizes. However, these were replaced with the visual mapping approach, as we found that it helped to cue users about the approximate position and function of each organ when deciding on clinical relevance. Our collaborators also found that using similar layouts for the dose-cluster encoding and additive effects view reduced cognitive load and made the system more visually consistent.

## 5.3 Outcome Plot

To support validation and iterative model improvement, it was important to show how outcomes vary within each cluster. This is important when ensuring, for example, that the cluster with the highest doses is actually capturing the high risk patients. To do this, we provide two types of encodings that show patient outcomes for each cluster: a temporal view of symptom ratings for the clusters, and a statistical bar-chart view showing the results of the likelihood ratio tests performed on each cluster for the outcome of interest.

Our temporal view uses a novel encoding (Fig. 3) to encode the trajectory of the symptom of interest across the entire treatment period for the patient clusters. This encoding has two components: a symbol grid, and a simple line chart. To reduce the complexity of the encoding, we first group the symptom ratings and treatment dates into bins (we selected five). In the symbol grid, we divide the patients into those in the selected cluster, and those not in the



**Figure 4:** Per-organ dose distribution for a selected cluster. Color gradients shows within-cluster distributions. (Left) A tooltip shows the full dose-volume histogram for a brushed organ. Dotted area shows the value (V55) currently being shown in the heat map.

selected cluster (out of cluster). For each patient, we calculate the highest rating for the symptom within the treatment dates before aggregating by cluster. We then calculate the percentage of patients from the selected cluster that fall in each rating + date bin. These percentages are encoded as circles on a grid, where the x-axis shows each date bin, and the y-axis encodes the symptom ratings. Size encodes the percentage of patients. Values for the in-cluster patients are shown as a saturated marker, while the out-of-cluster patients are shown as a black border marker. By comparing the markers, we can approximate the odds ratio of a patient within the selected cluster having a symptom of a given severity at each time point.

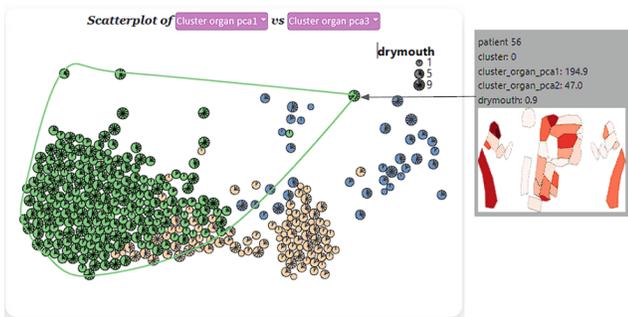
In addition to the symbol grid, we overlay a line chart that shows the mean symptom value over time for each cluster. The line charts use cluster colors. A cluster chart can be clicked to select that cluster for more details.

The statistical bar chart view encodes the results of the LRT test (Sec. 4.4). This view is used for assessing how well a model performs while accounting for the specific outcomes and confounders. Cluster-outcome relationships that are statistically significant ( $p < .05$ ) are shown using their categorical cluster color, while relationships that are not ( $p > .05$ ) are shown in gray. The selected cluster for the interface is highlighted using a bold black border between the bars for that cluster.

## 5.4 Cluster Dose-Distribution Plots

Once a reasonable set of cluster features has been identified, our first set of tasks involves investigating the dose distribution within each cluster (A2 T5-6). This is useful for identifying when the clusters are separating out patients with higher dose to other organs that were not included in the cluster inputs. To do this, we calculate the quartile ranges of a user-selected dose value within each cluster, for each organ. These values are then shown as a gradient heatmap using our 2-dimensional organ diagram using a sequential red color scheme (Sec. 5.1), where the innermost color represents the top quantile (80%) and the outer color represents the lower quantile (20%), allowing us to visualize the inter-organ dose distribution for each organ.

Interactions allow directly adding or removing organs from the cluster queue, as well as selecting a cluster to be used for brushing



**Figure 5:** Stylized scatterplot. Patients are represented by a custom glyph that encodes the outcome of interest (late drymouth ratings) as marks extending radially. Markers are colored by cluster membership, and a contour is shown around the currently selected cluster. A tooltip (left) shows a heatmap of the dose applied.

in other views. This facilitates the investigation of other aspects of the cluster in more detail.

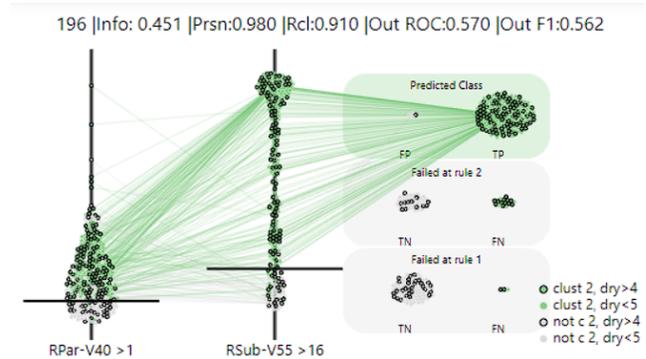
To anchor the visual heatmap in the clinician’s knowledge, we add a tooltip for each organ that can show the dose-volume histogram for each quantile for the selected organ and cluster (Fig. 4). This allows for a more detailed view of the entire histogram, while highlighting the relationship between the novel heatmap, and the standard dose-volume histogram that clinicians are familiar with.

### 5.5 Scatterplot

To visualize the distribution of patients across each cluster, we include a modified scatterplot panel that shows a 2-dimensional plot of the patients across two interactively-selected dimensions (Fig. 5). By default, we show the first two principal components of the features used to cluster the patients, but allow choosing to alternatively view higher order principal components, the principal components of the symptoms, or individual clinical or symptom ratings. Because we found that avoiding visual occlusion was more important than a high-fidelity projection, we use a force directed layout to remove overlap between glyphs.

Each patient in the scatterplot is encoded with a custom glyph that encodes its cluster membership, and the rating for the symptom of interest between 0 and 10. Each circular glyph is encoded with ticks that extend in 32.7-degree intervals in a clockwise radial pattern, where the number of ticks corresponds to the symptom rating. Thus, a full “pinwheel” glyph represents a patient with a symptom rating of 10, while an empty circle represents a patient that does not experience the symptom. Because symptom ratings use discrete ordinal (integer) values, we can encode the exact ratings. We additionally scale the size of the glyph based on the symptom rating to support visual identification of small or high dose values.

Finally, we color code the glyphs based on their cluster membership. The selected cluster is brushed by giving the corresponding glyphs a higher opacity, and drawing a contour around the convex hull of the cluster in the scatterplot. By hovering the mouse over a patient glyph, the user can view a tooltip showing a plot of the given patient’s received dose, and ratings for all symptoms over time. The dose to each organ is encoded for each patient using the organ diagram heatmap (Sec. 5.4).



**Figure 6:** Ruleset encoding in the rule mining view. A swarm plot of the patients is shown for the feature used in each rule, with the first and most informative rule on the right. A horizontal line shows the cutoff thresholds used in the rule. Patients that pass a rule are then plotted in a swarm plot in the next rule on the right. The section on the right shows rule patients failed at, with patients that pass all rules at the top (green area). Patients in each section are divided to show the False Positives or False Negatives at each level. Lines connect markers for a patient across each sub-plot.

Previous designs used alternative projection methods with alternative projections and glyph encodings. However, we found that allowing inspection of individuals was more important than preserving location with perfect fidelity. In contrast, T-SNE avoided occlusion, but tended to produce visual clusters that did not correspond to the desired clusters. For glyph design, we considered alternative shapes (e.g. diamonds or circles) for different levels of severity. However, collaborators found the use of color and shape confusing, while the use of ticks + size was better received and we were able to identify the patient of most interest (very high and very low severity) fairly easily for further inspection.

### 5.6 Rule Builder

Once our clusters are built, one of our goals is to explain the clusters in terms that are familiar to clinicians. To accomplish this, we used a constrained rule mining algorithm (Sec. 4.4) to produce a set of dose thresholds such that the group of patients that meet these thresholds approximates the selected cluster. This approach was chosen as clinicians often work with dose thresholds when choosing treatment plans.

When a cluster of interest is selected, our algorithm finds a list of rule-sets that optimize the mutual information between the patients and the cluster of interest. We then generate a plot for each ruleset, and show the top rules in a list to the user. We also show the number of predicted positives, information gain, precision, recall, and f1 for predicting the true class above each plot.

Our novel rule encoding is based on a mixture of swarm plots and parallel coordinate plots that are modified to show the progressive filtering of each ruleset (Fig. 6). We encode each feature (e.g. V50 to the tongue) along the x-axis. We then map the y-axis to the dose value in grays. Patients are plotted along the y-axis based on their value for the given dose feature in the x-axis, and adjusted using a force-directed layout to avoid overlap. A horizontal line is then drawn at the threshold of the rule for the feature on each step of the

x-axis. Patient marks are color-coded based on the selected cluster, while patients not in the selected cluster are gray.

To show the effect of additional rules, the features along the x-axis are ordered from left to right by the maximum information gain for its corresponding rule. In the first feature, we show all patients in the cohort. For additional features, we filter out all patients that do not satisfy rules from all previous features. The rightmost side of the encoding shows the patient groups stratified along the y-axis based on when they were filtered out of the ruleset. The set of patients that satisfy all rules is grouped at the top, while the set of patients that do not satisfy the first rule is grouped at the bottom. We further separate the final group by those in the true class (target cluster) and those not in the true class, allowing us to visualize the false positives and false negatives for each rule.

To provide a visual cue for how the rules are filtering the cohort along the x-axis, we provide lines that connect the undistorted locations of patients between axes, equivalent to a parallel coordinate plot with filtering. Once a patient is filtered out, we draw a line from the corresponding rule to the group on the right side. To prevent overlap, we only show the lines for the patients within one stratum at time, which is changed by brushing a patient in the given strata. By default, we brush the group of patients that satisfy all rules (predicted positives).

## 6 Evaluation

The first and foremost value of DASS comes from its unique functionality and its ability to support clinical model development, which we illustrate via two case studies. These case studies, presented here in abbreviated form, illustrate the process of creating models for practical use, based on real clinical data. The case studies were performed via Zoom meetings with desktop sharing, with one of the data scientists piloting DASS and the group using the think-aloud methodology with note-taking. We furthermore collected and report qualitative feedback from clinical collaborators during these case studies.

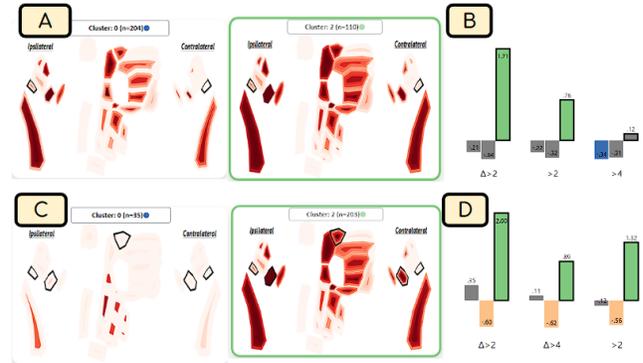
As further evidence of the DASS functional value, we provide in the supplemental materials a quantitative evaluation of clusters generated with DASS against baseline ML clusters. The DASS clusters improve performance for drymouth, choking, and swallowing issues. Finally, with an eye towards the generalizability of DASS to other modeling problems, we collected additional feedback where eight data scientists rated the usefulness and usability of DASS.

Since the interactive model-building components are directly targeted at modelers, An additional quantitative comparison of our clusters against baseline ML clusters generated without DASS can be found in the supplemental materials.

Our dataset consists of 349 patients treated with radiation therapy for oropharyngeal cancer. These models have been generated with the help of DASS by four data scientists in our group over several months of remote collaboration. The models have shown improvements over baseline models, and have been favorably evaluated by three clinical oncologists.

### 6.1 Case Study 1

Our group was interested in identifying patients at high risk of developing drymouth at 6 months after treatment, a common side



**Figure 7:** Case 1. (A) Low- and high-dose clusters using starting features. The low dose cluster includes several organs with high variance in the dose distribution. (B) Initial model performance. (C) Low- and high-dose clusters using the final model. Low dose cluster has a much lower variance with only a few sets of outliers. (D) Final model performance measures. High-risk cluster is correlated with drymouth with a higher odds ratio than the initial clusters.

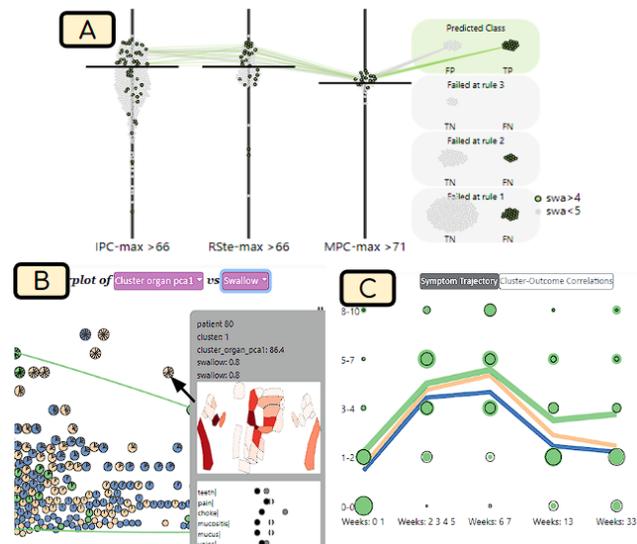
effect in HNC patients. In particular, the clinician analysts in the group wished to model the relationship between drymouth and the radiation dose applied to the salivary glands. The medical literature had established a few dose guidelines for parotid glands, but not for other salivary glands.

The model building process started by setting the parameters in the DASS control panel. Based on results from earlier work [WHL\*19], the group set the initial clustering features to be V40-V55 doses to the ipsilateral and contralateral Parotid glands. Three clusters based on a Gaussian mixture model were investigated. Inspecting the initial clusters in the outcome plot, the analysts noticed that, as expected, there was a higher rate of drymouth in the highest dose cluster (Cluster 2 in Fig. Sec. 7), although the correlation was not significant for the desired threshold of  $> 5$ . Moving to the dose distribution plot, the group noted that the low and medium dose clusters tended to have a high-variation in the dose to certain organs, as indicated by the dark red inner contours and light outer contours to several organs (Fig. 7-A), suggesting that the model parameters did not differentiate the low dose patients well. Moving to the additive effects view, the model was iteratively adjusted to include the submandibular glands and soft palate, with a larger dose window (V30-V55). After updating the model, the group noticed the clusters in the scatterplot panel achieved much better separation in the data (Fig. 7-D) compared to using just the parotid glands (Fig. 7-B). Returning to the dose cluster plots, the group also verified that the low dose cluster had a lower overall variance in the doses (Fig. 7-C).

Once the group achieved a set of features, the analysts aimed to verify the validity of the resulting model. Looking at the outcome panel, they noticed that while the high dose cluster was a strong predictor of drymouth, the low dose cluster had a high odds ratio. Moving back to the scatterplot, and with the help of the oncologists, they inspected the patients in this low dose group, and noticed an interesting pattern: a number of patients had very high symptom ratings, and confirmed that none of their organs received notably

high doses. Pivoting to the temporal outcome panel, the analysts further noted that this low-dose group had the highest incidence of severe drymouth at the start of treatment. After further discussion with the clinical collaborators, the group concluded that existing treatment plans try to minimize dose to the parotid glands, but not the submandibular glands, so the dose tends to be much lower in severe cases. The team theorized that there is likely a minor, but not full compensatory effect of the contralateral salivary glands when one set of salivary glands fails that should be explored later when investigating dose guidelines.

## 6.2 Case Study 2



**Figure 8:** Case 2. (A) Rule mining results for predicting severe swallow dysfunction, which suggest using high doses to the pharyngeal constrictors. (B) Scatterplot of the first principle component of the cluster features vs swallow ratings. A tooltip highlights a case with severe swallowing in a low-dose cluster. (C) Outcome plot for the final clusters. High risk patients have similar ratings during treatment, but swallowing issues increase between 6 weeks and 6 months after treatment.

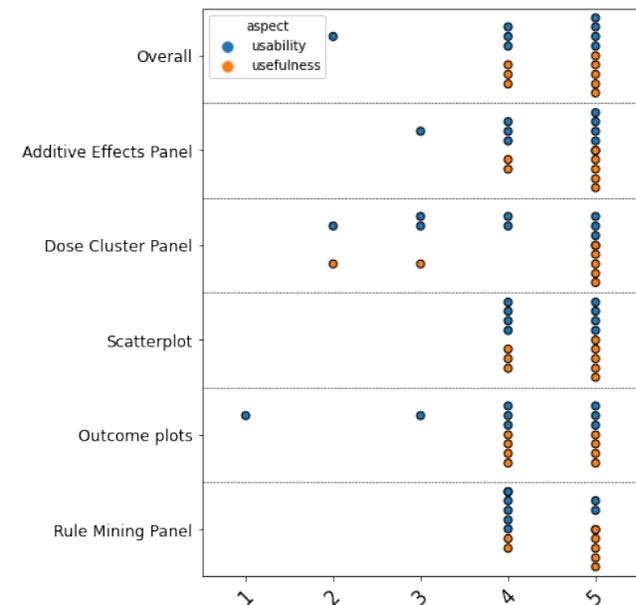
This second case study dealt with the identification of patients at high risk of swallowing dysfunction, which is a less common outcome that is theorized to be related to damage to muscles in the mouth and throat. Swallowing disorders are also related to patients that require a feeding tube and weight loss, and thus it is an important outcome to avoid. High-risk patients can also be assigned prehabilitative therapy such as swallowing exercises as well.

To help identify a set of starting organs, the analysts inspected the rule mining view and set the desired outcome to be severe late swallowing using all available features (Fig. 8-A). By looking at the resulting rules, the group was able to identify the organs and dose features that best predicted severe swallowing, which allowed selecting a set of starting features for the cluster. Among the best splits were high dose depths (V55-V70) to the superior, medial, and inferior pharyngeal constrictors, which are key muscles used in swallowing, which were chosen as a starting point for the clusters. After running the clustering, the analysts inspected the out-

come view and noticed that the initial clustering parameters were effectively separating the high-risk patients: this highest dose cluster had a significantly higher odds ratio of severe late swallowing (2.56) than other clusters (Fig. 8-C). Inspecting the cluster dose distribution view, it was noted that this high-dose cluster was noticeably smaller ( $n = 35$ ) than the drymouth cluster and that the high-dose cluster tended to consistently have a much higher V55 to the IPC than other clusters.

Moving to the scatterplot, the analysts changed the dimensions to show the first principle component of the dose and swallowing ratings, which allowed identifying all patients with high swallow dysfunction that were not in the cluster (Fig. 8-B). Using the tooltip, the group found some of these patients had high doses to the base of the tongue and upper larynx. The analysts then added the supra-glottic larynx to the clustering parameters in hopes of capturing this group. The group then moved to the additive effects view, iteratively changed the dose window to include only the V55-V65, and added the esophagus, which is another major muscle used for swallowing in the base of the throat. After finalizing the parameter set, the analysts inspected the rule view to find the features that best distinguished the high-risk cluster. This high-risk cluster was easily distinguished using the V55 to the Inferior Pharyngeal Constrictor. Our clinical collaborators noted that all the pharyngeal constrictor muscles are located close together, and there exist guidelines for the dose for all of these muscles. Thus, a high IPC dose is likely a predictor of a high dose to all related organs. Additionally, the group discussed the fact that the dose threshold for swallowing was higher than drymouth, which may indicate that muscles are less sensitive to radiation relative to salivary glands.

## 6.3 General Usefulness and Usability Feedback



**Figure 9:** General DASS usability and usefulness.

In addition to the case studies, which illustrate the DASS unique functionality, we collected qualitative and quantitative feedback

from both collaborators and from modelers not affiliated with the project. All collaborators appreciated the functionality provided by DASS, and are in the process of publishing the resulting clinical models. Regarding the spatial cluster panel, our clinical collaborators found it intuitive and useful for inspecting dose distributions of organs of interest. Feedback on the rule mining algorithm was also positive, with oncologists remarking that it was “very useful”, as it could “translate our results into practical applications”. A data mining expert responded similarly to the additive effects panel, saying that it was a “nice, very nice way to explore the parameter space”.

Additionally, we asked, via an anonymous online questionnaire, three senior data scientists in the group, who were not directly involved in the DASS design but participated in walkthroughs of the system, and five junior data scientists, who were not affiliated with the project to rate the usefulness and usability of the whole system and of each component of the system on a Likert scale from 1 to 5. We specifically sought feedback from data scientists, with an eye towards generalizability, as modelers are the intended users of the interactive model-building components of the system. Results are shown in Fig. 9.

Feedback was very positive, with most ratings between 4-5, in particular for usefulness. Ratings for usability were slightly lower, as expected: some of the group experts clarified the group’s narration during the model-building process was extremely useful, and they wished for visual help buttons replicating that experience on demand. Ratings from the junior data scientists not affiliated with the project were occasionally lower, in particular for the composite outcome marker plots and the dose cluster panel. Based on the qualitative feedback, the difference in these cases was directly related to the visual scaffolding and domain expertise which collaborators benefited from, as these plots were based on methods used in RT planning and clinical biostatistics.

## 7 Discussion and Conclusion

Our design relies on three main principles for improving model development: 1) information scent to guide model development (A1); 2) visual scaffolding to support bridging the information gap between what domain experts commonly deal with and what is needed to reason about the data (A2); 3) model explanations aimed at translating our novel approach to the types of simpler “models” use in practice (A3). Our case studies show how the system was effectively used to develop explainable models that outperformed our previous attempts at developing clinical models.

Below, we distill the design lessons gathered from this project when dealing with visual steering and explainable AI problems in collaboration with domain experts.

**L1. Explanation Scaffolding:** We extend the concept of visual scaffolding – gradually building to more complex visualizations from a more familiar one – to that of XAI-style model explanations. Specifically, we argue that model explanations should aim to translate more complex models into those that mimic how users commonly deal with the data. In our case, we used constrained rule mining in conjunction with visualizing intra-cluster dose distributions using a visual scaffolding approach. Other systems have used regression models which are common in biostatistics. However clin-

icians do not often reason about such models directly, so they are less useful in clinical practice.

**L2. Keep Model Goals Flexible:** When developing models, data scientists may work solely to optimize the performance in terms of easily measured outcomes [MWM\*19], which leads to issues during collaboration with model end-users [ZMW20b]. In practice, there is often a misalignment between what can easily be measured, and what makes a model useful in practice. In developing our models, we found that it was important to allow users to investigate a mixture of outcomes, in addition to qualitative factors such as model plausibility and complexity, which need to be leveraged against each other when deciding on the final model.

**L3. Encourage Skepticism:** One motivation in the design of our system was a recurring problem of designing models that performed well, whereas further probing revealed internal logic that appeared to be the result of biases and spurious correlations in the data. Despite this, our models were often received without skepticism when these issues were not brought up. This issue with over-trusting erroneous explanations has been suggested in early empirical studies [KNJ\*20, XSHF19]. The communication gap between model builders and experts may result in dramatically over-trusting the models for both parties as they may be unable to identify issues in the models on their own. When dealing with XAI, designers should focus on promoting skepticism about the models by highlighting potential issues in the models, such as outliers and confounders, which can help highlight previously unknown issues in the models.

The main limitation of our system is the reliance on visualizations that require familiarity as well as knowledge of the underlying models and data, which is made possible by the long-term nature of our collaboration. While we use domain-specific designs for our visual scaffolding approach and model designs, the design philosophy can be generalized to other problems involving spatial data where model outputs can be translated into discrete groups, such as clustering and decision trees. In terms of scalability, our system requires 5-15 seconds to update new results for each cluster, depending on the number of clusters and rules mining settings. Scaling to larger datasets may increase the required time, although this is still significantly faster than alternatives that do not use interactive steering. Visualization of individual patients in the Scatterplot and Rule view may also be difficult with very large cohorts.

In conclusion, we have presented an ML and visual steering system for clinical oncology symptom modeling with spatial data. We described the co-design of a clinical visual-steering system, and demonstrated its ability to support the creation of interpretable ML models for stratifying patients. Additionally, we presented a set of lessons learned for model co-development and model explanations for a hybrid, machine expert and human expert problem. We hope that these findings will help future designers create better, and more trustworthy models in high-stakes settings.

**Acknowledgements** Our work is supported by NIH awards NCI-R01-CA258827 and NLM-R01-LM012527, and NSF awards CDSE-1854815 and CNS-1828265.

## References

[ASO\*16] ANGULO D. A., SCHNEIDER C., OLIVER J. H., CHARPAK N., HERNANDEZ J. T.: A multi-faceted visual analytics tool for ex-

- ploratory analysis of human brain and function datasets. *Front. neuroinformatics* 10 (2016), 36. doi:10.3389/fninf.2016.00036. 2
- [BSM\*15] BERNARD J., SESSLER D., MAY T., SCHLOMM T., PEHRKE D., KOHLHAMMER J.: A visual-interactive system for prostate cancer cohort analysis. *Comp. Graph. and App.* 35, 3 (2015), 44–55. doi:10.1109/MCG.2015.49. 2
- [CD18] CAVALLO M., DEMIRALP Ç.: Clustrophile 2: Guided visual clustering analysis. *Trans. Vis. Comp. Graph.* 25, 1 (2018), 267–276. doi:10.1109/TVCG.2018.2864477. 2
- [CLD\*21] CHENG F., LIU D., DU F., LIN Y., ZYTEK A., LI H., QU H., VEERAMACHANENI K.: Vbridge: Connecting the dots between features and data to explain healthcare models. *Trans. Vis. Comp. Graph.* 28, 1 (2021), 378–388. doi:10.1109/TVCG.2021.3114836. 3
- [CLL\*21] CHOI J., LEE S.-E., LEE Y., CHO E., CHANG S., JEONG W.-K.: Dxplore: A unified visualization framework for interactive dendritic spine analysis using 3d morphological features. *Trans. Vis. Comp. Graph.* (2021), 1–1. doi:10.1109/TVCG.2021.3116656. 2
- [CMQ20] CHENG F., MING Y., QU H.: Dece: Decision explorer with counterfactual explanations for machine learning models. *Trans. Vis. Comp. Graph.* 27, 2 (2020), 1438–1447. doi:10.1109/TVCG.2020.3030342. 3
- [DvVH\*19] DINGEN D., VAN’T VEER M., HOUTHUIZEN P., MESTROM E. H. J., KORSTEN E. H., BOUWMAN A. R., VAN WIJK J.: Regressionexplorer: Interactive exploration of logistic regression models with subgroup analysis. *Trans. Vis. Comp. Graph.* 25, 1 (2019), 246–255. doi:10.1109/TVCG.2018.2865043. 3
- [ELB\*91] EMAMI B., LYMAN J., BROWN A., COLA L., GOITEIN M., MUNZENRIDER J., SHANK B., SOLIN L., WESSON M.: Tolerance of normal tissue to therapeutic irradiation. *Int. Jour. Rad. Onco. Bio. Phys.* 21, 1 (1991), 109–122. doi:10.1016/0360-3016(91)90171-y. 2
- [ENBD08] EL NAQA I., BRADLEY J. D., DEASY J. O.: Nonlinear kernel-based approaches for predicting normal tissue toxicities. In *IEEE Int. Conf. on Mach. Learn. and App.* (2008), pp. 539–544. doi:10.1109/ICMLA.2008.126. 2
- [FGM\*20] FURMANOVÁ K., GROSSMANN N., MUREN L. P., CASARES-MAGAZ O., MOISEENKO V., EINCK J. P., GRÖLLER M. E., RAIDOU R. G.: Vapor: visual analytics for the exploration of pelvic organ variability in radiotherapy. *Comp. & Graph.* 91 (2020), 25–38. doi:10.1016/j.cag.2020.07.001. 2
- [FH15] FEHRENBACH M. J., HERRING S. W.: *Illustrated Anatomy of the Head and Neck*, vol. 5. Elsevier Health Sciences, 2015. 5
- [FMCM\*21] FURMANOVÁ K., MUREN L. P., CASARES-MAGAZ O., MOISEENKO V., EINCK J. P., PILSKOG S., RAIDOU R. G.: Previs: Predictive visual analytics of anatomical variability for radiotherapy decision support. *Comp. & Graph.* 97 (2021), 126–138. doi:10.1016/j.cag.2021.04.010. 2
- [FNB\*21] FLORICEL C., NIPU N., BIGGS M., WENTZEL A., CANAHUATE G., VAN DIJK L., MOHAMED A., FULLER C. D., MARAI G. E.: Thalix: Human-machine analysis of longitudinal symptoms in cancer therapy. *Trans. Vis. Comp. Graph.* 28, 1 (2021), 151–161. doi:10.1109/TVCG.2021.3114810. 2
- [FYTL18] FALK M., YNNERMAN A., TREANOR D., LUNDSTRÖM C.: Interactive visualization of 3d histopathology in native resolution. *Trans. Vis. Comp. Graph.* 25, 1 (2018), 1008–1017. doi:10.1109/TVCG.2018.2864816. 2
- [GCMM\*19] GROSSMANN N., CASARES-MAGAZ O., MUREN L. P., MOISEENKO V., EINCK J. P., GRÖLLER M. E., RAIDOU R. G.: Pelvis runner: Visualizing pelvic organ variability in a cohort of radiotherapy patients. In *Eurographics Work. on Vis. Comp. for Bio. and Med.* (2019), pp. 69–78. doi:10.2312/vcbm.20191233. 2
- [GGC\*17] GLUECK M., GVOZDIK A., CHEVALIER F., KHAN A., BRUDNO M., WIGDOR D.: Phenosticks: Cross-sectional cohort phenotype comparison visualizations. *Trans. Vis. Comp. Graph.* 23, 1 (2017), 191–200. doi:10.1109/TVCG.2016.2598469. 2
- [GNDV\*17] GLUECK M., NAEINI M. P., DOSHI-VELEZ F., CHEVALIER F., KHAN A., WIGDOR D., BRUDNO M.: Phenolines: Phenotype comparison visualizations for disease subtyping via topic models. *Trans. Vis. Comp. Graph.* 24, 1 (2017), 371–381. doi:10.1109/TVCG.2017.2745118. 2
- [GRW\*00] GRESH D. L., ROGOWITZ B. E., WINSLOW R. L., SCOLLAN D. F., YUNG C. K.: Weave: A system for visually linking 3-d and statistical visualizations, applied to cardiac simulation and measurement data. In *IEEE Vis.* (2000), pp. 489–492. doi:10.1109/VISUAL.2000.885739. 2
- [GXZ\*17] GUO S., XU K., ZHAO R., GOTZ D., ZHA H., CAO N.: Eventthread: Visual summarization and stage analysis of event sequence data. *Trans. Vis. Comp. Graph.* 24, 1 (2017), 56–65. doi:10.1109/TVCG.2017.2745320. 2
- [JBB\*08] JAINEK W. M., BORN S., BARTZ D., STRASSER W., FISCHER J.: Illustrative hybrid visualization and exploration of anatomical and functional brain data. *Comp. Graph. For.* 27, 3 (2008), 855–862. doi:https://doi.org/10.1111/j.1467-8659.2008.01217.x. 2
- [JBF\*20] JÖNSSON D., BERGSTRÖM A., FORSELL C., SIMON R., ENGSTRÖM M., WALTER S., YNNERMAN A., HOTZ I.: Visualneuro: A hypothesis formation and reasoning application for multi-variate brain cohort study data. *Comp. Graph. For.* 39, 6 (2020), 392–407. doi:10.1111/cgf.14045. 2
- [JKW\*22] JESSUP J., KRUEGER R., WARCHOL S., HOFFER J., MUHLICH J., RITCH C. C., GAGLIA G., COY S., CHEN Y.-A., LIN J.-R., SANTAGATA S., SORGER P. K., PFISTER H.: Scope2screen: Focus+context techniques for pathology tumor assessment in multivariate image data. *Trans. Vis. Comp. Graph.* 28, 1 (2022), 259–269. doi:10.1109/TVCG.2021.3114786. 2
- [KEV\*17] KWON B. C., EYSENBACH B., VERMA J., NG K., DE FILIPPI C., STEWART W. F., PERER A.: Clustervision: Visual supervision of unsupervised clustering. *Trans. Vis. Comp. Graph.* 24, 1 (2017), 142–151. doi:10.1109/TVCG.2017.2745085. 2
- [KK08] KONISHI S., KITAGAWA G.: *Information criteria and statistical modeling*. Springer, 2008. 4
- [KNJ\*20] KAUR H., NORI H., JENKINS S., CARUANA R., WALLACH H., WORTMAN VAUGHAN J.: Interpreting interpretability: understanding data scientists’ use of interpretability tools for machine learning. In *Proc. ACM on Human-Comp. Int.* (2020), pp. 1–14. doi:10.1145/3313831.3376219. 10
- [KOJL\*14] KLEMM P., OELTZE-JAFRA S., LAWONN K., HEGENSCHIED K., VÖLZKE H., PREIM B.: Interactive visual analysis of image-centric cohort study data. *Trans. Vis. Comp. Graph.* 20, 12 (2014), 1673–1682. doi:10.1109/TVCG.2014.2346591. 2
- [KPS16] KRAUSE J., PERER A., STAVROPOULOS H.: Supporting iterative cohort construction with visual temporal queries. *Trans. Vis. Comp. Graph.* 22, 1 (2016), 91–100. doi:10.1109/TVCG.2015.2467622. 2
- [LDVdL\*08] LANGENDIJK J. A., DOORNAERT P., VERDONCK-DE LEEUW I. M., LEEMANS C. R., AARONSON N. K., SLOTMAN B. J.: Impact of late treatment-related toxicity on quality of life among patients with head and neck cancer treated with radiotherapy. *Jour. Clin. Onco.* 26, 22 (2008), 3770–3776. doi:10.1200/JCO.2007.14.6647. 2
- [LWE\*20] LUCIANI T., WENTZEL A., ELGOHARI B., ELHALAWANI H., MOHAMED A., CANAHUATE G., VOCK D. M., FULLER C. D., MARAI G. E.: A spatial neighborhood methodology for computing and analyzing lymph node carcinoma similarity in precision medicine. *Journal of biomedical informatics* 112 (2020), 100067. doi:10.1016/j.yjbix.2020.100067. 2
- [MAM20] MUSTAQEEM A., ANWAR S. M., MAJID M.: A modular cluster based collaborative recommender system for cardiac patients. *Art. Intel. Med.* 102 (2020), 101761. doi:10.1016/j.artmed.2019.101761. 2

- [Mar17] MARAI G. E.: Activity-centered domain characterization for problem-driven scientific visualization. *Trans. Vis. Comp. Graph.* 24, 1 (2017), 913–922. doi:10.1109/TVCG.2017.2744459. 3
- [MDM\*15] MALIK S., DU F., MONROE M., ONUKWUGHA E., PLAISANT C., SHNEIDERMAN B.: Cohort comparison of event sequences with balanced integration of visual analytics and statistics. In *20th Int Conf on Int. User Int.* (2015), pp. 38–49. doi:10.1145/2678025.2701407. 2
- [MK22] MARWAHA J. S., KVEDAR J. C.: Crossing the chasm from model performance to clinical impact: the need to improve implementation and evaluation of ai, 2022. doi:10.1038/s41746-022-00572-2. 2
- [MKKW12] MARLIN B. M., KALE D. C., KHEMANI R. G., WETZEL R. C.: Unsupervised pattern discovery in electronic health care data using probabilistic clustering models. In *2nd ACM SIGHIT Int. Health Info. Symp.* (2012), Association for Computing Machinery, p. 389–398. doi:10.1145/2110363.2110408. 2
- [MMB\*18] MARAI G. E., MA C., BURKS A. T., PELLOLIO F., CANAHUATE G., VOCK D. M., MOHAMED A. S., FULLER C. D.: Precision risk analysis of cancer therapy with interactive nomograms and survival plots. *Trans. Vis. Comp. Graph.* 25, 4 (2018), 1732–1745. doi:10.1109/TVCG.2018.2817557. 3
- [MP13] MÜHLBACHER T., PIRINGER H.: A partition-based framework for building and validating regression models. *Trans. Vis. Comp. Graph.* 19, 12 (2013), 1962–1971. doi:10.1109/TVCG.2013.125. 3
- [MPL\*18] MA C., PELLOLIO F., LLANO D. A., STEBBINGS K. A., KENYON R. V., MARAI G. E.: Rembrain: Exploring dynamic biospatial networks with mosaic matrices and mirror glyphs. *Elec. Imag.* (2018), 1–13. doi:10.2352/J.ImagingSci.Technol.2017.61.6.060404. 2
- [MQB18] MING Y., QU H., BERTINI E.: Rulematrix: Visualizing and understanding classifiers with rules. *Trans. Vis. Comp. Graph.* 25, 1 (2018), 342–352. doi:10.1109/TVCG.2018.2864812. 3
- [MV15] METSALU T., VILO J.: Clustvis: a web tool for visualizing clustering of multivariate data using principal component analysis and heatmap. *Nucleic acids res.* 43, W1 (2015), W566–W570. doi:10.1093/nar/gkv468. 2
- [MWM\*19] MAO Y., WANG D., MULLER M., VARSHNEY K. R., BALDINI I., DUGAN C., MOJSILOVIĆ A.: How data scientists work together with domain experts in scientific collaborations: To find the right answer or to ask the right question? *Proc. ACM Human-Comp. Int.* 3 (2019), 1–23. doi:10.1145/3361118. 10
- [MXC\*19] MING Y., XU P., CHENG F., QU H., REN L.: Protosteer: Steering deep sequence model with prototypes. *trans. vis. comp. graph.* 26, 1 (2019), 238–248. doi:10.1109/TVCG.2019.2934267. 3
- [ODH\*07] OELTZE S., DOLEISCH H., HAUSER H., MUIGG P., PREIM B.: Interactive visual analysis of perfusion data. *Trans. Vis. Comp. Graph.* 13, 6 (2007), 1392–1399. doi:10.1109/TVCG.2007.70569. 2
- [OHS\*16] O’SULLIVAN B., HUANG S. H., SU J., GARDEN A. S., STURGIS E. M., DAHLSTROM K., LEE N., RIAZ N., PEI X., KOYFMAN S. A., ET AL.: Development and validation of a staging system for hpv-related oropharyngeal cancer by the international collaboration on oropharyngeal cancer network for staging (icon-s): a multi-centre cohort study. *The Lancet Oncol.* 17, 4 (2016), 440–451. doi:10.1016/S1470-2045(15)00560-4. 2
- [Raf95] RAFTERY A. E.: Bayesian model selection in social research. *Sociological methodology* (1995), 111–163. doi:10.2307/271063. 4
- [RCMA\*18] RAIDOU R. G., CASARES-MAGAZ O., AMIRKHAPOV A., MOISENKO V., MUREN L. P., EINCK J. P., VILANOVA A., GRÖLLER M. E.: Bladder runner: Visual analytics for the exploration of rt-induced bladder toxicity in a cohort study. In *Comp. Graph. For.* (2018), vol. 37, Wiley Online Library, pp. 205–216. doi:10.1111/cgf.13413. 2
- [RCMM\*16] RAIDOU R., CASARES-MAGAZ O., MUREN L., VAN DER HEIDE U., RØRVIK J., BREEUWER M., VILANOVA A.: Visual analysis of tumor control models for prediction of radiotherapy response. *Comp. Graph. For.* 35, 3 (2016), 231–240. doi:https://doi.org/10.1111/cgf.12899. 3
- [RMC\*07] ROSENTHAL D. I., MENDOZA T. R., CHAMBERS M. S., ASPER J. A., GNING I., KIES M. S., WEBER R. S., LEWIN J. S., GARDEN A. S., ANG K. K., ET AL.: Measuring head and neck cancer symptom burden: the development and validation of the md anderson symptom inventory, head and neck module. *Head & Neck* 29, 10 (2007), 923–931. doi:10.1002/hed.20602. 3
- [RPHL14] RISTOVSKI G., PREUSSER T., HAHN H. K., LINSEN L.: Uncertainty in medical visualization: Towards a taxonomy. *Comp. & Graph.* 39 (2014), 60–73. doi:10.1016/j.cag.2013.10.015. 2
- [RvdHD\*15] RAIDOU R., VAN DER HEIDE U., DINH C., GHOBADI G., KALLEHAUGE J., BREEUWER M., VILANOVA A.: Visual analytics for the exploration of tumor tissue characterization. *Comp. Graph. For.* 34, 3 (2015), 11–20. doi:https://doi.org/10.1111/cgf.12613. 3
- [SMS\*22] STREEB D., METZ Y., SCHLEGEL U., SCHNEIDER B., EL-ASSADY M., NETH H., CHEN M., KEIM D. A.: Task-based visual interactive modeling: Decision trees and rule-based classifiers. *Trans. Vis. Comp. Graph.* 28, 9 (2022), 3307–3323. doi:10.1109/TVCG.2020.3045560. 3
- [STA\*22] SRABANTI S., TRAN M., ACHIM V., FULLER D., CANAHUATE G., MIRANDA F., MARAI G. E.: A tale of two centers: Visual exploration of health disparities in cancer care. In *Pac. Vis. Symp. (PacificVis)* (2022), IEEE, pp. 101–110. doi:10.1109/pacificvis53943.2022.00019. 2
- [TM03a] TEOH S. T., MA K.-L.: Paintingclass: interactive construction, visualization and exploration of decision trees. In *Int. Conf. on Knowl. disc. and data min.* (2003), pp. 667–672. doi:10.1145/956750.956837. 3
- [TM03b] TEOH S. T., MA K.-L.: Starclass: Interactive visual classification using star coordinates. In *Int. Conf. on Data Min.* (2003), SIAM, pp. 178–185. doi:10.1137/1.9781611972733.16. 3
- [vdEvW11] VAN DEN ELZEN S., VAN WIJK J. J.: Baobabview: Interactive construction and analysis of decision trees. In *Conf. on Vis. Anal. Sci. and Tech. (VAST)* (2011), pp. 151–160. doi:10.1109/VAST.2011.6102453. 3
- [WCVD\*20] WENTZEL A., CANAHUATE G., VAN DIJK L. V., MOHAMED A. S., FULLER C. D., MARAI G. E.: Explainable spatial clustering: Leveraging spatial data in radiation oncology. In *Vis. Conf. (short paper)* (2020), IEEE, pp. 281–285. doi:10.1109/VIS47514.2020.00063. 2, 5
- [WHC\*22] WANG Q., HUANG K., CHANDAK P., ZITNIK M., GEHLENBORG N.: Extending the nested model for user-centric xai: A design study on gnn-based drug repurposing. *Trans. Vis. Comp. Graph.* 29, 1 (2022), 1266–1276. doi:10.1109/TVCG.2022.3209435. 3
- [WHL\*19] WENTZEL A., HANULA P., LUCIANI T., ELGOHARI B., ELHALAWANI H., CANAHUATE G., VOCK D., FULLER C. D., MARAI G. E.: Cohort-based t-ssim visual computing for radiation therapy prediction and exploration. *Trans. Vis. Comp. Graph.* 26, 1 (2019), 949–959. doi:10.1109/TVCG.2019.2934546. 1, 2, 5, 8
- [WHvD\*20] WENTZEL A., HANULA P., VAN DIJK L. V., ELGOHARI B., MOHAMED A. S., CARDENAS C. E., FULLER C. D., VOCK D. M., CANAHUATE G., MARAI G. E.: Precision toxicity correlates of tumor spatial proximity to organs at risk in cancer patients receiving intensity-modulated radiotherapy. *Radiotherapy and Onc.* 148 (2020), 245–251. doi:10.1016/j.radonc.2020.05.023. 1, 2, 3
- [WKN\*12] WARCHOL S., KRUEGER R., NIRMAL A. J., GAGLIA G., JESSUP J., RITCH C. C., HOFFER J., MUHLICH J., BURGER M. L., JACKS T., SANTAGATA S., SORGER P. K., PFISTER H.: Viscinity: Visual spatial neighborhood analysis for multiplexed tissue imaging data. *Trans. Vis. Comp. Graph.* (1912), 1–11. doi:10.1109/TVCG.2022.3209378. 2

- [WMH\*21] WANG Q., MAZOR T., HARBIG T. A., CERAMI E., GEHLENBORG N.: Threadstates: State-based visual analysis of disease progression. *Trans. Vis. Comp. Graph.* 28, 1 (2021), 238–247. doi:10.1109/TVCG.2021.3114840. 2
- [WVW\*04] WOLF I., VETTER M., WEGNER I., NOLDEN M., BOTTGER T., ET AL.: The medical imaging interaction toolkit (MITK): a toolkit facilitating the creation of interactive software by extending VTK and ITK. In *Med. Imag.: Vis., Image-Guided Proc., and Disp.* (2004), pp. 16–28. doi:10.1117/12.535112. 2
- [XSFM11] XU W., SMITH A. M., FAEDER J. R., MARAI G. E.: Rulebender: a visual interface for rule-based modeling. *Bioinformatics* 27, 12 (2011), 1721–1722. doi:10.2312/evs.20201067. 3
- [XSHF19] XIONG C., SHAPIRO J., HULLMAN J., FRANCONERI S.: Illusion of causality in visualized data. *Trans. Vis. Comp. Graph.* 26, 1 (2019), 853–862. doi:10.1109/TVCG.2019.2934399. 10
- [YNB21] YUAN J., NOV O., BERTINI E.: An exploration and validation of visual factors in understanding classification rule sets. In *Vis. Conf. (short paper)* (2021), IEEE, pp. 6–10. doi:10.1109/VIS49827.2021.9623303. 3
- [ZGP15] ZHANG Z., GOTZ D., PERER A.: Iterative cohort analysis and exploration. *Info. Vis.* 14, 4 (2015), 289–307. doi:10.1177/1473871614526077. 2
- [ZHT\*13] ZHANG L., HUB M., THIEKE C., FLOCA R. O., KARGER C. P.: A method to visualize the uncertainty of the prediction of radiobiological models. *Phys. Med.* 29, 5 (2013), 556–561. doi:10.1016/j.ejmp.2012.11.004. 3
- [ZMP\*21] ZHANG T., MCCOY T. H., PERLIS R. H., DOSHI-VELEZ F., GLASSMAN E.: Interactive cohort analysis and hypothesis discovery by exploring temporal patterns in population-level health records. In *IEEE Vis. Ana. in Heal. (VAHC)* (2021), pp. 14–18. doi:10.1109/VAHC53616.2021.00007. 2
- [ZMW\*20a] ZEBRALLA V., MÜLLER J., WALD T., BOEHM, ET AL.: Obtaining patient-reported outcomes electronically with “oncofunction” in head and neck cancer patients during aftercare. *Front. in Onco.* 10 (2020), 2502. doi:10.3389/fonc.2020.549915. 2
- [ZMW20b] ZHANG A. X., MULLER M., WANG D.: How do data science workers collaborate? roles, workflows, and tools. *Proc. ACM on Human-Comp. Int.* 4 (2020), 1–23. doi:10.1145/3392826. 10

## A. Workflow and Background

Fig A1. Physician Workflow. The main items of interest for this project are to help establish a stratification of patients' risk of certain side-effects from their radiation plan (Side-Effect red box), which can then be used to identify which patients require additional preventative treatment, as well as help identify dose thresholds (Dose thresholds red box) that can be input as soft-constraints during treatment planning.

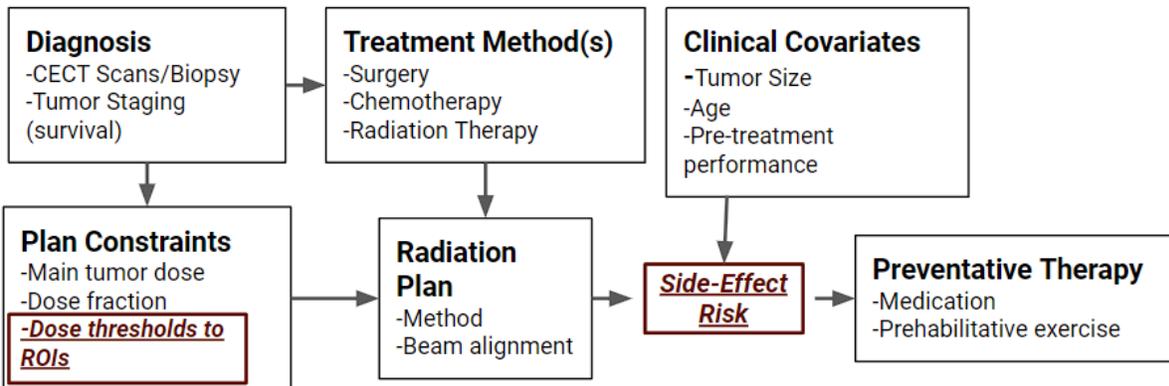


Fig A2. Full 2D Mapping of Regions of Interest. Regions represent approximate locations of important anatomical structures along the center, and sides of the head. The relative size of smaller regions are exaggerated to improve visibility.

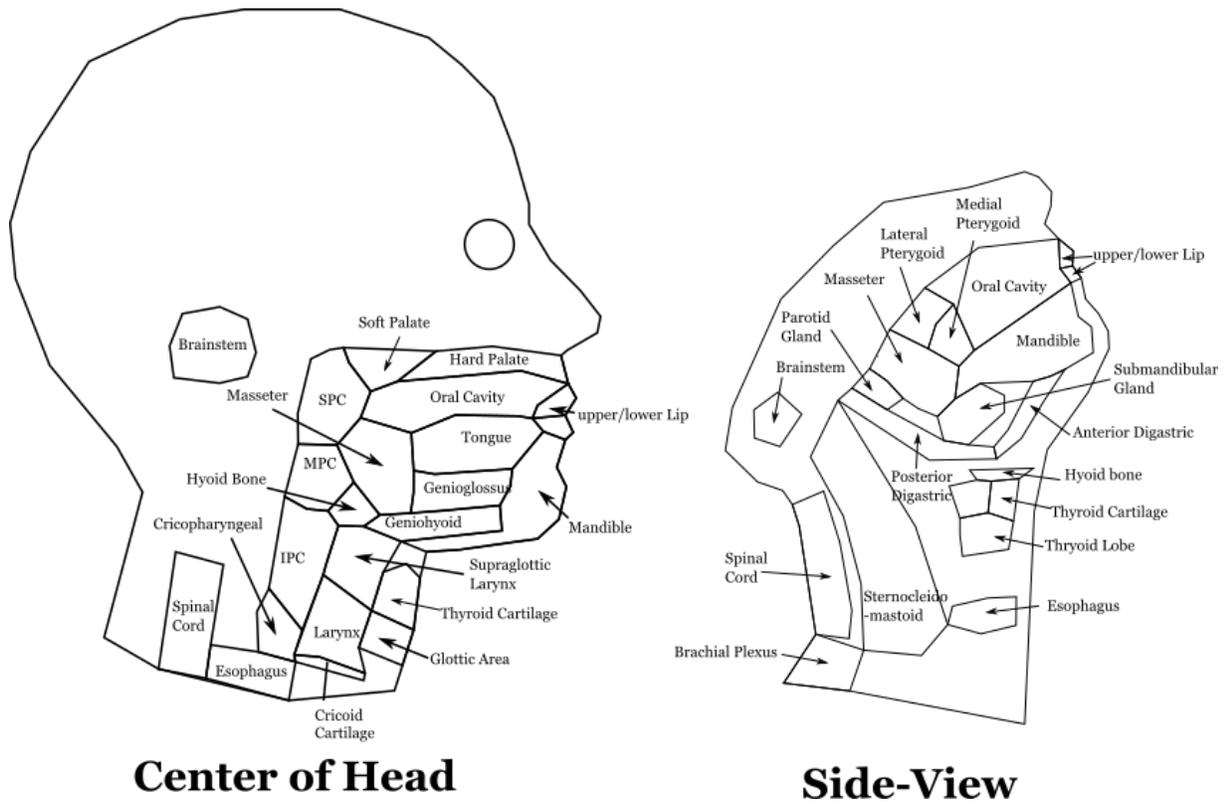
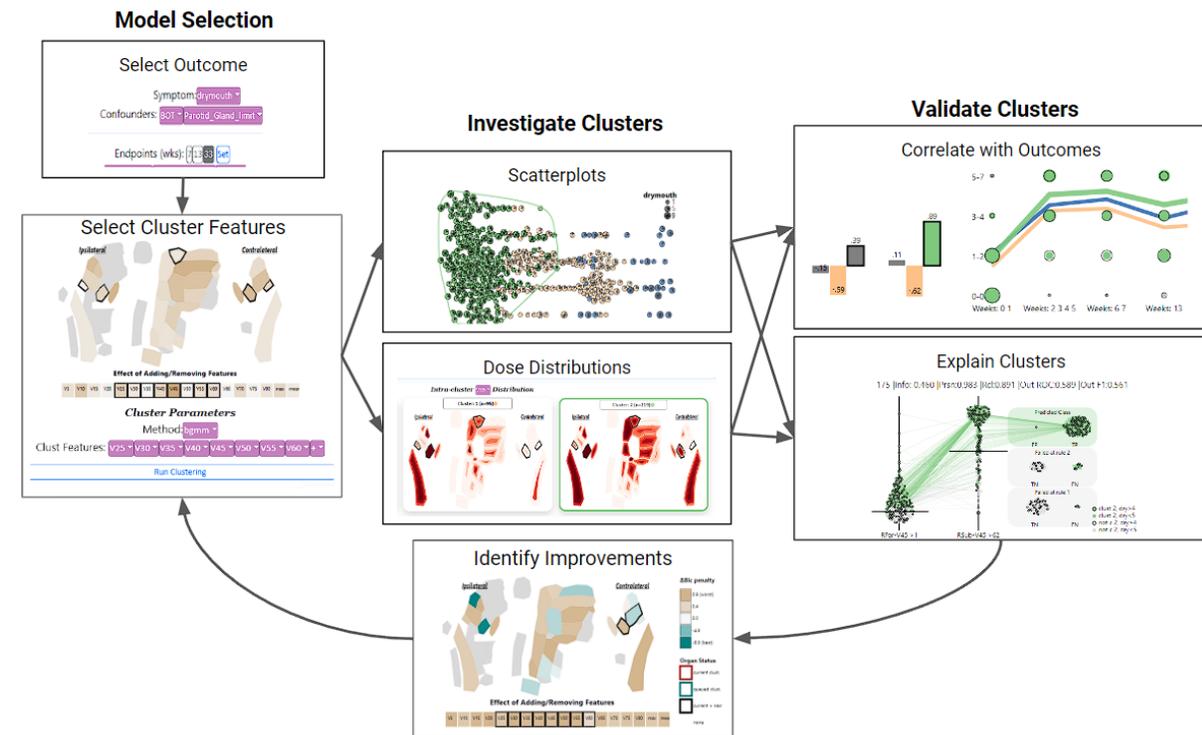


Fig A3. DASS Workflow. First, a desired outcome is selected, along with the initial clustering parameters, which can be drawn from prior literature or informed from information metrics in the explanation view. After clusters are generated, the inter-cluster distributions can be investigated using the dose distribution and configurable scatterplot views. Cluster performance can then be validated by investigating the inter-cluster symptom trajectories and correlations with outcomes. A rule-based classifier can also be used to produce explanations for the high or low-risk clusters based on dose thresholds. Once clusters are investigated, the Additive Effects panel can identify potential changes to the clustering parameters that could improve the model performance.



## B. Prototypes and Designs

This section details some of the early prototypes and approaches that led up to our final design of DASS, which covers both algorithmic and visualization attempts.

Fig B1. Early prototype of the visual scaffolding interface used to inspect the results of clusters generated using grid-search or physician-selected features without visual steering. Each row shows the distribution of mean-dose to each ROI within each cluster. Bar charts on the left show odds-ratios of different symptoms. ROIs used in the cluster use a saturated red color scheme, while ROIs not used in the clusters use a desaturated grayscale color scheme. This was our original approach that inspired the design of MOTIV, as we felt that our clusters were not producing outcomes that achieved both good performance and explanations that aligned with clinical knowledge.





Fig B4. View that shows the effect of each cluster on a machine learning model for predicting different presence of a symptom at 6 months at different thresholds. A user selected model is trained to predict each outcome using clinical variables. Additional models are built that also include cluster labels (all or individual cluster labels), and cross-validation performance is measured. Outcomes are grouped by the metric used, and the dose threshold considered (3,5,7). This view was removed as our collaborators preferred to rely on traditional statistics with fewer thresholds when deciding on outcomes.

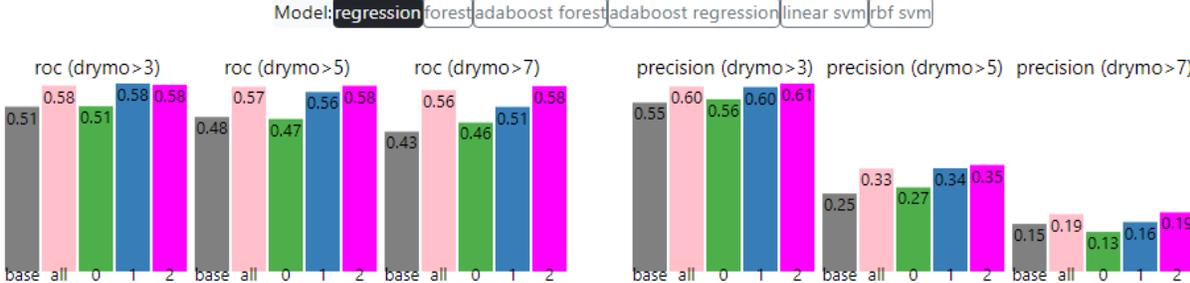


Fig B5. Early version of the interface using patient-specific view and star-charts to encode symptom ratings for patients and clusters. The patient view (top-right) used a single selected patient (133) within a selected cluster at the top. The additional patients on the left are the most similar within the same cluster, while the patients on the left show similar patients in different clusters, based on treatment categories and doses. Patients in the cluster are encoded with a user-selected dose value (e.g. V55) to each organ in red, while patients in other clusters are color-encoded with the difference in dose values from the selected patient using a blue-green color scale, where green indicates higher dose, and blue indicates lower doses. Star charts encode patient symptoms. This view was removed as we found that inspecting patients via a tooltip on the scatterplot was more efficient.

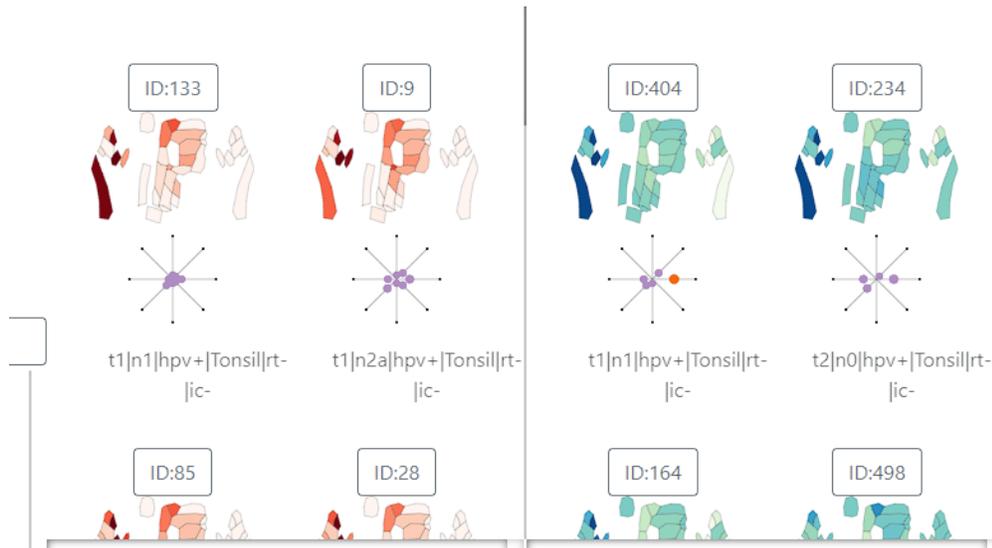


Fig B6. Alternative Variants of the patient encodings that use a color gradient to show the full dose-histogram within each organ. The outer layer represents the V5 while the innermost layer represents the V75 to the respective organ. Symptoms at different stages (baseline, during treatment, and post-treatment) are shown as bar charts using small multiples for a set of symptoms of interest. Patient color encodings for the counterfactual patients show absolute rather than relative doses.

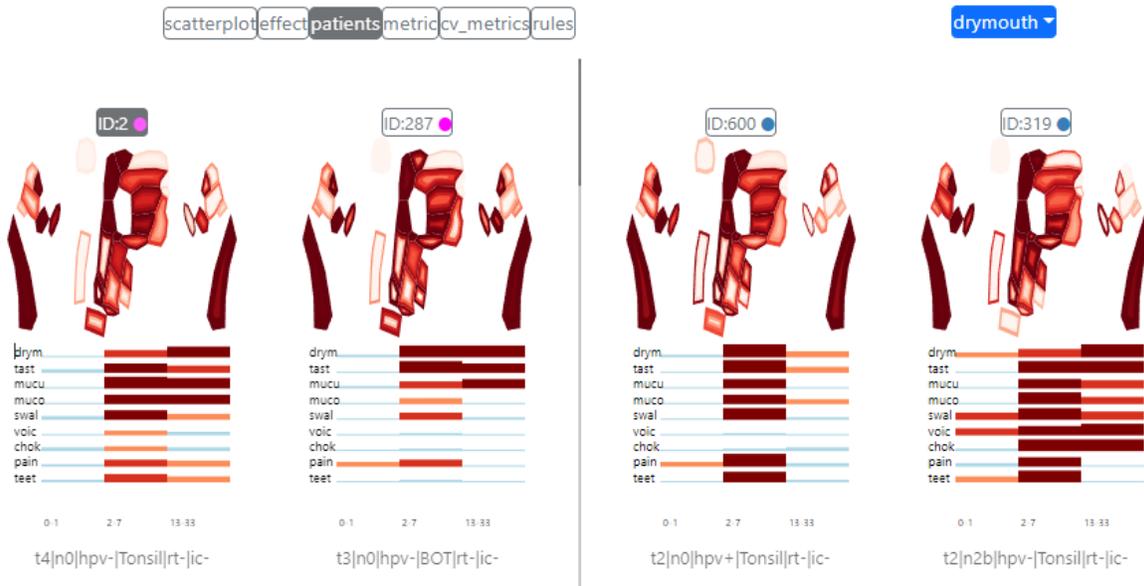


Fig B7. Alternative Version of the scatterplot that used shape instead of glyphs to encode symptom ratings. Circles, triangles, and diamonds encoded patients with symptom ratings of < 5, 5-7, or >7, respectively. Shape-encodings were changed to glyphs as our collaborators felt that the encodings were unintuitive.



Fig B8. early version of the cluster view. The view on the right encodes cluster doses as in the final version. A star chart on the center-right encodes symptom ratings within the cluster, where the dot encodes median rating and the bold line indicates the 25-75% confidence interval. Circle size encodes odds-ratio for the cluster to experience severe symptoms. Finally, histograms on the right encode the distribution of clinical features within each cluster. The symptom charts were replaced with our current symptom plots as we wanted to focus on comparing trajectories between clusters. Clinical feature plots were removed to save space, as we did not use the information in the development of clusters outside of using them as covariates in our predictive models. Instead, extensive demographic distributions were calculated after clusters were finalized and thus did not need to be included in the interactive system.

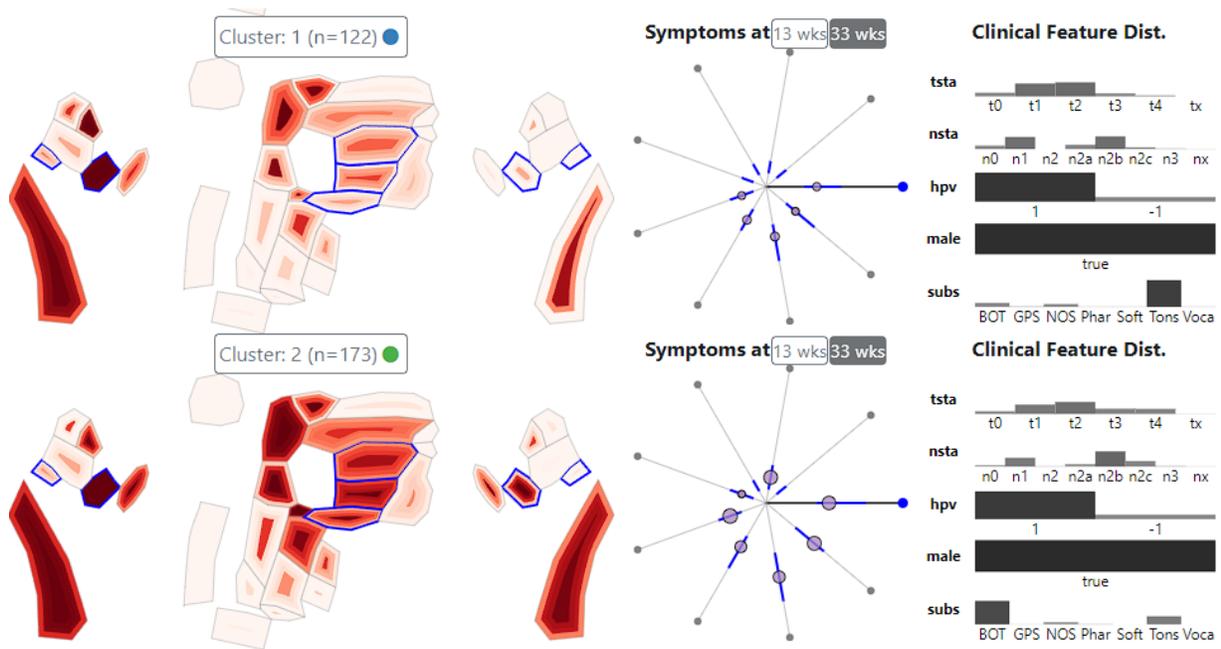
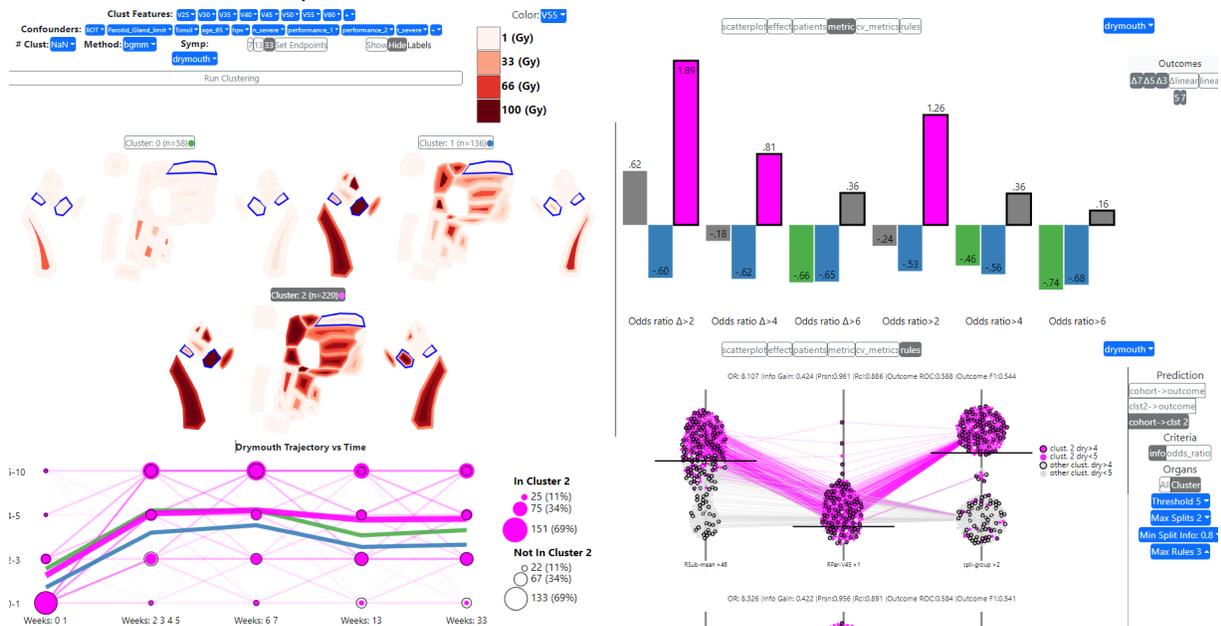
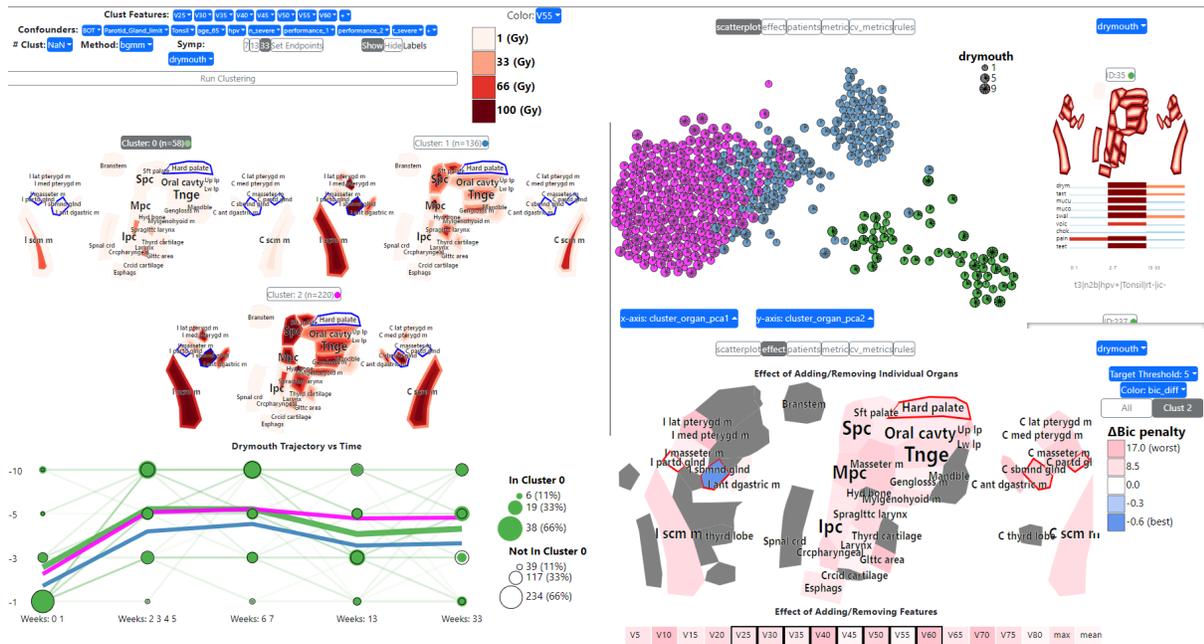


Fig B9. Early design of the layout of DASS in different configurations. The original layout was divided into 5 views: a control panel (top-left), the cluster view (center left), the symptom trajectory view (bottom-left), and two windows on the right that allow the user to toggle between views. Several adjustments to the layout and color schemes were made to de-clutter the interface based on feedback from collaborators. We also incorporated control-panels for each view into titles for each plot in order to make them more readable.





## C. Cluster Quantitative Analysis

We include here a comparison of the performance of the clusters contained using our system to baseline clusters that are obtained via methods without visual steering. In this experiment, we use a baseline method of generating clusters to predict patient outcomes without use of the DASS interface. We then compare those results to clusters made from the DASS interface using two kinds of starting parameters: a default set determined by treatment guidelines, and the parameters used in the baseline clusters. Clusters are evaluated using a simple ablation study where a regression model is built to predict the desired outcome with and without the cluster label to test the effect of the model on cross validation model performance. We perform this experiment on three outcomes: late severe (rating >5) drymouth, swallow dysfunction, and choking.

To generate baseline clusters, we tested the results of clusters generated through a combination of input features and organs with 2-4 clusters, using a bayesian mixture model. Because we require that the organs in the model be linked to the outcome, we compiled lists of potential organ combinations for each outcome based on clinical literature. For drymouth, we consider the set of organs taken from NTCP models built for predicting xerostomia [1], and for swallow and choke, we use organs taken from NTCP models built for dysphagia [2]. In addition, we included sets of organs gathered from clinician feedback on the organs they believe would be related to each.

To compare our models, we used two variants of clusters that are obtained from DASS for each outcome. As a baseline, we found clusters following the regular approach, where we used a general starting point of the most important organ of interest. In addition, we consider models

that are obtained from Dass by using the model parameters in the baseline models found using grid-search. Simplified clusters for all cases are generated as described in the main text. The parameters of the final clusters are listed in Table C1.

For quantitative analysis, we convert our clusters into a clinical stratification during cross-validation as follows: first, we rank each cluster based on the number of patients that experience the given outcome in the training data split and assign risk to patients in the test data based on the rank of their clusters, where the highest-risk cluster is given a risk score of 1, the second highest-risk cluster is given a score of  $(\text{number of clusters} - 2) / (\text{number of clusters} - 1)$ , etc, and the lowest risk is given a score of 0. For the simplified cluster, we always assign a risk of 1 to the high-dose cluster and 0 otherwise. We then calculate the cross-validation mean Area-under the curve score (AUC), and the Mathews correlation coefficient (MCC) for predicting each outcome for each cluster model. In addition, we calculate the results from using only the highest-dose (HD) cluster for each stratification, as this is a likely use case for clinical applications when deciding whether to prescribe preventative treatment. Thus, we generate AUC and MCC scores for 3 different cases for each clustering.

Results for AUC are shown in Figure C1. Drymouth clusters obtained using DASS, the main symptom of interest that we used to build the model, perform better for all measures compared to grid-search. Choke clusters perform better in general, although interestingly, grid-search simplified explanations tend to outperform other rule-based simplified explanations. For swallow, the second Dass clusters outperform the gridsearch cluster for the stratification but not the rule-based clusters.

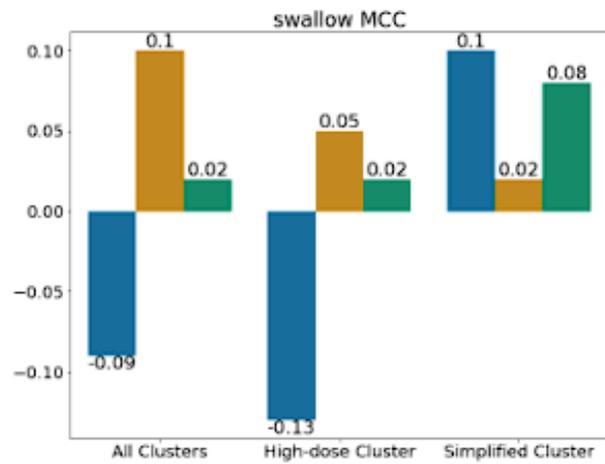
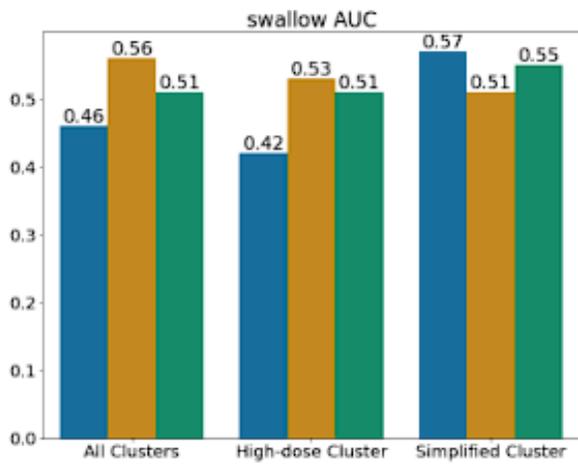
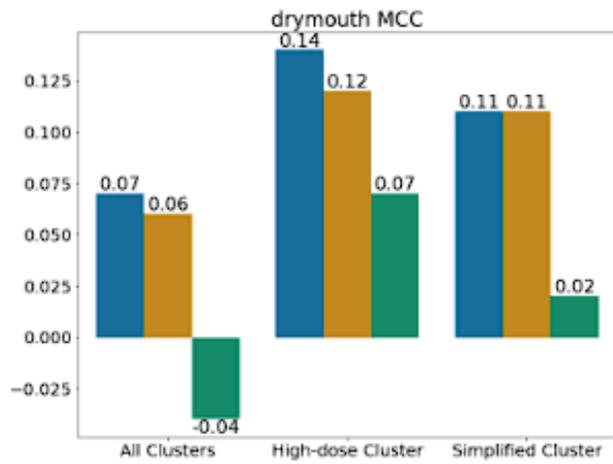
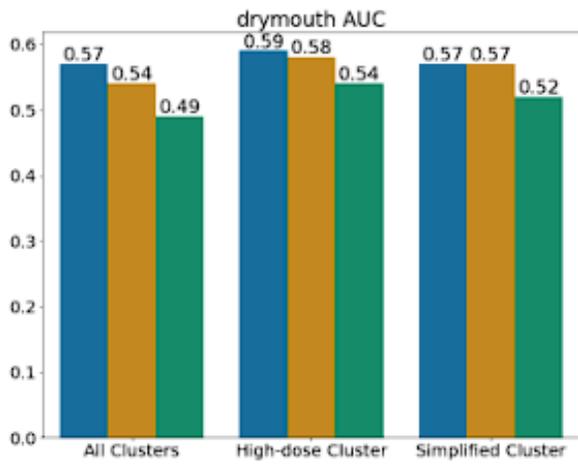
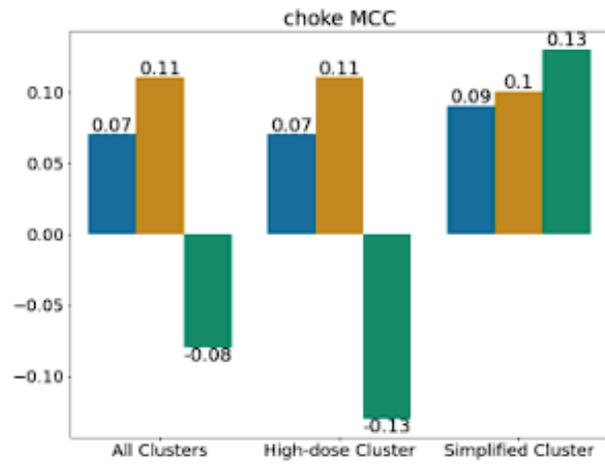
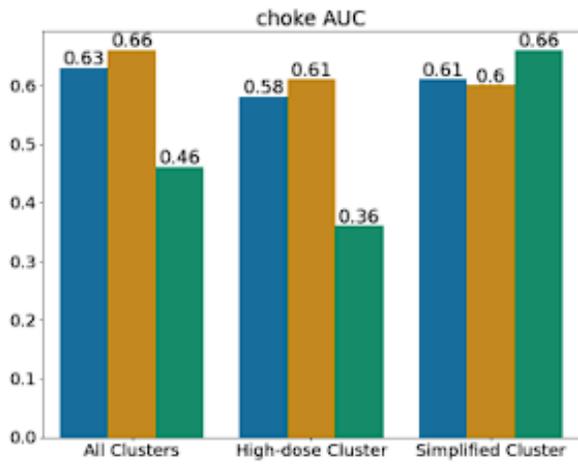
We note that while our swallow and choke clusters have mixed results, we can reasonably beat the baseline clusters for several measures without the need for extensive literature review to identify the optimal organs to consider when testing multiple clusters, in addition to the general benefits of interactive changing of the clusters to balance between other considerations such as the maximum number of clusters or comparison of performance across multiple thresholds. Drymouth results tended to be the most consistent, most likely due to the higher rates of severe drymouth, while severe swallow and choke are relatively rare, making the cluster predictions difficult. In addition, the organs linked to drymouth tend to be salivary glands rather than muscles, and are likely more sensitive to lower doses.

Table C1. Table of features used in each cluster for our analysis. The prefixes 'Lt' and "Rt" refer to the ipsilateral (same side as main tumor) and contralateral (opposite side as main tumor) organ, respectively. Legend: Dass: Models built using all clusters using DASS with a default starting point; Alt Dass: Models built using DASS with the baseline clusters at a starting point. Baseline: Models built using clusters obtained through standard grid search.

Symptom	organs	N Clusters	Organs	Dose Features	Simple Rule
<b>Drymouth</b>	Dass	3	Rt_Parotid_Gland, Lt_Parotid_Gland, Rt_Submandibular_Gland, Lt_Submandibular_Gland, Hard_Palate	V25-V60	V55_Rt_Submandibular_Gland>23.0, V65_Lt_Submandibular_Gland>51.0, V10_Rt_Parotid_Gland>18.0
	Dass Alt	3	Rt_Parotid_Gland, Lt_Parotid_Gland, Rt_Submandibular_Gland, Lt_Submandibular_Gland, Soft_Palate	V30-V55	V55_Rt_Submandibular_Gland>23.0, V65_Lt_Submandibular_Gland>51.0, V10_Rt_Parotid_Gland>18.0
	Baseline	3	Lt_Parotid_Gland, Rt_Parotid_Gland, Lt_Submandibular_Gland, Rt_Submandibular_Gland, Soft_Palate, Upper_Lip, Lower_Lip, Extended_Oral_Cavity, Mylogeniohyoid_M	mean dose, max dose	V60_Mylogeniohyoid_M>32.0, V55_Extended_Oral_Cavity>28.0, V30_Rt_Submandibular_Gland>62.0, V65_Lt_Submandibular_Gland>1.0
<b>Choke</b>	Dass	3	IPC, MPC, Supraglottic_Larynx, Esophagus, Mylogeniohyoid_M	V30-V65	V60_Supraglottic_Larynx>40.0, V55_MPC>17.0, V55_Mylogeniohyoid_M>24.0, V15_Esophagus>22.0
	Dass Alt	3	IPC, MPC, Supraglottic_Larynx, Cricopharyngeal_Muscle, Mylogeniohyoid_M	V30, V40, V50, V60	mean_dose_IPC>49.0, mean_dose_MPC>54.0, mean_dose_Cricopharyngeal_Muscle>16.0
	Baseline	3	SPC, IPC, Supraglottic_Larynx, Rt_Parotid_Gland, Cricopharyngeal_Muscle	mean dose, max dose	V5_Rt_Parotid_Gland>72.0, V65_SPC>3.0, V55_Supraglottic_Larynx>46.0, max_dose_IPC>46.0
<b>Swallow</b>	Dass	3	IPC, MPC, SPC' Supraglottic_Larynx, Esophagus	V35-V70	V60_Supraglottic_Larynx>40.0, V65_SPC>3.0, V15_Esophagus>22.0,

					V55_MPC>17.0
	Dass Alt	2	IPC, SPC, Cricopharyngeal_Muscle, Rt_Parotid_Gland	V65, max dose	max_dose_Cricopharyn geal_Muscle>63.0
	Baseline	2	SPC, IPC, Supraglottic_Larynx, Rt_Parotid_Gland, Cricopharyngeal_Muscle	mean_dose , V40-V65	V5_Rt_Parotid_Gland>7 2.0, V65_SPC>3.0, V55_Supraglottic_Laryn x>46.0, max_dose_IPC>46.0

Figure C1. AUC and MCC scores for all model types. For cluster stratifications, Alt Dass clusters and high-dose clusters outperform the baseline models. Alt Dass performs better than Dass for Choke and Swallow, but not Drymouth.



**Legend**



DASS Clusters



DASS Clusters 2



Baseline

Table C2. Results from 5-fold cross-validation of different stratification models for predicting severe late symptoms. We consider both absolute ratings > 4 at 6 months, and increase in ratings > 4 from baseline at 6 months. Models include stratification with all clusters (stratification), stratification using only the highest dose clusters (High-dose), and the simplified high-risk rule explanations. Legend: Dass: Models built using all clusters using DASS with a default starting point; Alt Dass: Models built using DASS with the baseline clusters at a starting point. Baseline: Models built using clusters obtained through standard grid search.

Symptom	Change From Baseline	Model	ROC-AUC	MCC	Precision	Recall	F1
	FALSE	Dass Stratification	0.628	0.070	0.059	0.500	0.106
		Alt Dass Stratification	0.659	0.106	0.080	0.429	0.135
		Baseline Stratification	0.462	-0.082	0.026	0.357	0.048
		Dass High-dose	0.584	0.070	0.059	0.500	0.106
		Alt Dass High-dose	0.611	0.106	0.080	0.429	0.135
		Baseline High-dose	0.357	-0.126	0.000	0.000	0.000
		Dass Simplified	0.611	0.088	0.060	0.643	0.110
		Alt Dass Simplified	0.596	0.100	0.083	0.357	0.135
		Baseline Simplified	0.662	0.127	0.066	0.786	0.122
		Dass Stratification	0.652	0.065	0.051	0.500	0.092
		Alt Dass Stratification	0.714	0.131	0.080	0.500	0.138
		Baseline Stratification	0.469	-0.085	0.021	0.333	0.039

Choke

TRUE

		Dass High-dose	0.584	0.065	0.051	0.500	0.092
		Alt Dass High-dose	0.648	0.131	0.080	0.500	0.138
		Baseline High-dose	0.358	-0.116	0.000	0.000	0.000
		Dass Simplified	0.623	0.090	0.053	0.667	0.099
		AI Dass Simplified	0.627	0.122	0.083	0.417	0.139
		Baseline Simplified	0.685	0.135	0.060	0.833	0.112
		Dass Stratification	0.571	0.065	0.220	0.574	0.318
		Alt Dass Stratification	0.540	0.058	0.216	0.603	0.318
		Baseline Stratification	0.490	-0.037	0.179	0.426	0.252
		Dass High-dose	0.586	0.137	0.244	0.691	0.360
		Alt Dass High-dose	0.577	0.124	0.236	0.706	0.354
		Baseline High-dose	0.542	0.069	0.215	0.706	0.330
		Dass Simplified	0.568	0.108	0.241	0.574	0.339
		AI Dass Simplified	0.568	0.108	0.241	0.574	0.339
	FALSE	Baseline Simplified	0.515	0.025	0.208	0.382	0.269
		Dass Stratification	0.583	0.046	0.141	0.568	0.226
		Alt Dass Stratification	0.546	0.035	0.137	0.591	0.222
		Baseline Stratification	0.485	-0.059	0.105	0.386	0.165
		Dass High-dose	0.613	0.150	0.171	0.750	0.278
		Alt Dass High-dose	0.596	0.130	0.163	0.750	0.267
		Baseline High-dose	0.551	0.070	0.143	0.727	0.240
		Dass Simplified	0.573	0.097	0.160	0.591	0.252

Drymouth

TRUE

		AI Dass Simplified	0.573	0.097	0.160	0.591	0.252
		Baseline Simplified	0.516	0.022	0.136	0.386	0.201
		Dass Stratification	0.465	-0.091	0.140	0.393	0.206
		Alt Dass Stratification	0.560	0.104	0.247	0.328	0.282
		Baseline Stratification	0.513	0.024	0.180	0.770	0.292
		Dass High-dose	0.418	-0.125	0.122	0.311	0.175
		Alt Dass High-dose	0.533	0.053	0.190	0.689	0.298
		Baseline High-dose	0.512	0.021	0.179	0.770	0.291
		Dass Simplified	0.567	0.102	0.218	0.557	0.313
	FALSE	AI Dass Simplified	0.511	0.018	0.179	0.754	0.289
	FALSE	Baseline Simplified	0.550	0.075	0.205	0.557	0.300
		Dass Stratification	0.506	-0.021	0.116	0.465	0.186
		Alt Dass Stratification	0.538	0.083	0.173	0.326	0.226
		Baseline Stratification	0.523	0.040	0.130	0.837	0.225
		Dass High-dose	0.431	-0.092	0.090	0.326	0.141
		Alt Dass High-dose	0.510	0.014	0.127	0.651	0.212
		Baseline High-dose	0.494	-0.011	0.122	0.837	0.212
		Dass Simplified	0.577	0.101	0.160	0.581	0.251
		AI Dass Simplified	0.504	0.007	0.125	0.744	0.213
Swallow	TRUE	Baseline Simplified	0.560	0.079	0.151	0.581	0.239

Table C3. Distribution of outcomes in the cohort considered.

Symptom	Total >4 at 6M	Mean Rating	95% CI
Drymouth	153 (43.84%)	4.338109	0.4-9
Swallow	46 (12.18%)	2.137536	0-7
Choke	19 (5.44%)	1.111748	0-5

Table C4. Parameters used in the gridsearch for the baseline. All parameters were tested with a gaussian mixture models using 2-4 components

Dose Values (All Outcomes)
['mean_dose','V25','V30','V35','V40','V45','V50','V55','V60','V65','V70','V75','V80'], ['mean_dose','V25','V30','V35','V40','V45','V50'], ['mean_dose','V25','V30','V35','V40'], ['mean_dose','V30','V35','V40','V45','V50','V55','V60','V65','V70','V75','V80'], ['mean_dose','V40','V45','V50','V55','V60','V65','V70','V75','V80'], ['mean_dose','V55','V60','V65','V70','V75','V80'], ['mean_dose','V60','V65','V70','V75','V80'], ['mean_dose','V40','V45','V50','V55','V60','V65'], ['mean_dose','V40','V45','V50','V55','V60'], ['V25','V30','V35','V40','V45','V50','V55','V60','V65','V70','V75','V80'], ['V25','V30','V35','V40','V45','V50'], ['V25','V30','V35','V40'], ['V30','V35','V40','V45','V50','V55','V60','V65','V70','V75','V80'], ['V40','V45','V50','V55','V60','V65','V70','V75','V80'], ['V55','V60','V65','V70','V75','V80'], ['V60','V65','V70','V75','V80'], ['V40','V45','V50','V55','V60','V65'], ['V40','V45','V50','V55','V60'], ['mean_dose'], ['max_dose'], ['max_dose','mean_dose']
Organs (Drymouth)

['Tongue', 'Mylogeniohyoid_M', 'Genioglossus_M', 'Rt_Parotid_Gland', 'Lt_Parotid_Gland', 'Rt_Submandibular_Gland', 'Lt_Submandibular_Gland', 'Soft_Palate', 'Extended_Oral_Cavity', 'Supraglottic_Larynx', 'Larynx']
['Lt_Parotid_Gland', 'Rt_Parotid_Gland', 'Lt_Submandibular_Gland', 'Rt_Submandibular_Gland', 'Soft_Palate', 'Upper_Lip', 'Lower_Lip', 'Extended_Oral_Cavity', 'Mylogeniohyoid_M']
Organs (Swallow)
['SPC', 'IPC', 'Supraglottic_Larynx', 'Rt_Parotid_Gland', 'Cricopharyngeal_Muscle']
['IPC', 'MPC', 'SPC', 'Mylogeniohyoid_M', 'Tongue']
['SPC', 'IPC', 'MPC', 'Supraglottic_Larynx', 'Rt_Parotid_Gland', 'Cricopharyngeal_Muscle']
Organs (Choke)
['SPC', 'IPC', 'Supraglottic_Larynx', 'Rt_Parotid_Gland', 'Cricopharyngeal_Muscle']
['Rt_Masseter_M', 'Lt_Masseter_M', 'Rt_Medial_Pterygoid_M', 'Lt_Medial_Pterygoid_M', 'Rt_Lateral_Pterygoid_M', 'Lt_Lateral_Pterygoid_M']
['Rt_Masseter_M', 'Lt_Masseter_M', 'Rt_Medial_Pterygoid_M', 'Lt_Medial_Pterygoid_M', 'Rt_Lateral_Pterygoid_M', 'Lt_Lateral_Pterygoid_M', 'Supraglottic_Larynx', 'Larynx', 'Glottic_Area', 'Thyroid_cartilage', 'Cricopharyngeal_Muscle', 'Cricoid_cartilage', 'Esophagus']

[1] Beetz, Ivo, et al. "NTCP models for patient-rated xerostomia and sticky saliva after treatment with intensity modulated radiotherapy for head and neck cancer: the role of dosimetric and clinical factors." *Radiotherapy and Oncology* 105.1 (2012): 101-106.

[2] Kierkels, Roel GJ, et al. "Multivariable normal tissue complication probability model-based treatment plan optimization for grade 2–4 dysphagia and tube feeding dependence in head and neck radiotherapy." *Radiotherapy and Oncology* 121.3 (2016): 374-380.

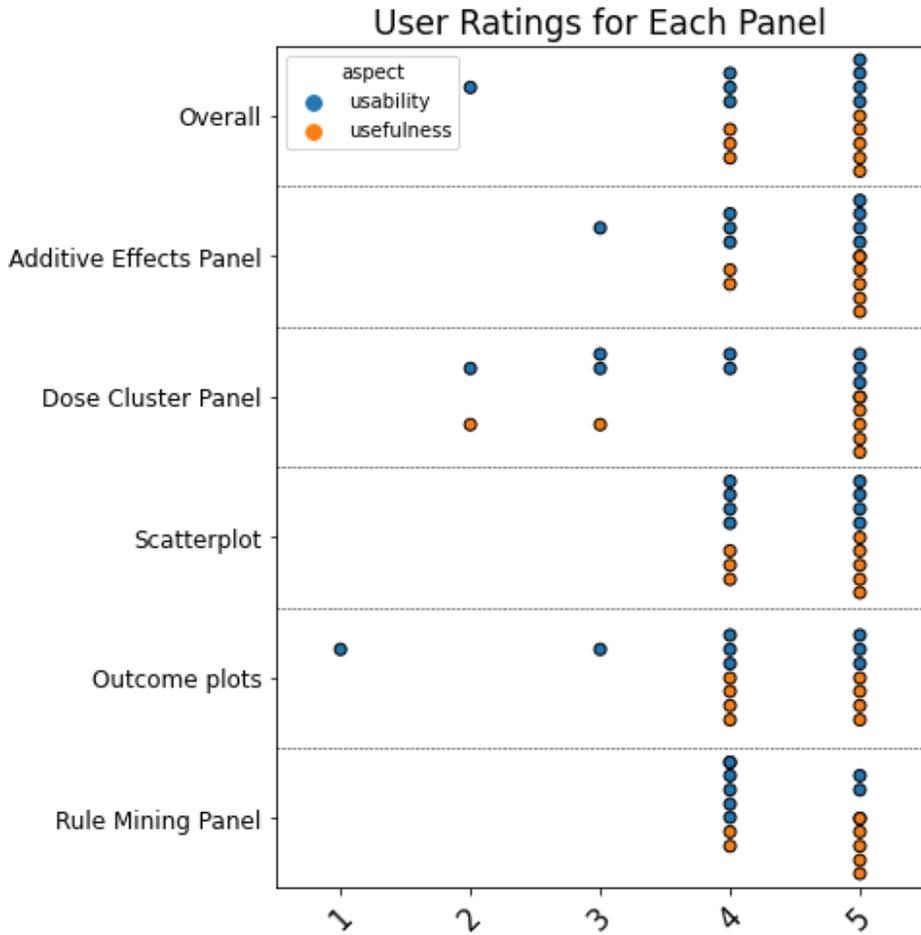
## D. Qualitative Analysis Results

Table D1. Results From the usability evaluation questionnaire. Ratings are given on a range of 1 (not useful) to 5 (very useful).

	Additive Effects	Dose Clusters	Scatterplot	Outcome Plots	Rule Mining	Overall
--	------------------	---------------	-------------	---------------	-------------	---------

	Explanations											
	Useful	Usable	Useful	Usable	Useful	Usable	Useful	Usable	Useful	Usable	Useful	Usable
P1	4	4	5	4	5	4	5	4	5	4	5	4
P2	5	5	5	5	5	5	4	5	5	4	5	5
P3	5	5	5	5	5	5	4	3	4	4	5	5
P4	5	5	5	5	5	5	4	4	4	4	5	5
P5	5	4	5	3	5	5	5	5	5	4	4	4
P6	5	4	3	3	4	4	5	4	5	5	4	4
P7	4	3	2	2	4	4	4	1	5	4	4	2
P8	5	5	5	4	4	4	5	5	5	5	5	5
Average	4.75	4.38	4.38	3.88	4.63	4.50	4.50	3.88	4.75	4.25	4.63	4.25

Figure D1. Plot of user responses to the usability questionnaire.



## Open ended comments

**P6:** The Dose Cluster Panel a bit complex at the beginning. But became clearer later on. Use cases are well explained in the end.

For the additive effects panel, it's easy to see the effect of adding other organs but I believe there is a factor of domain knowledge missing. For instance, adding an organ may look like a good idea but it may not be relevant to the symptom being investigated (based on AW's comments during presentation).

**P7:** The Dose Cluster panel encoding confused me because the distribution was mapped to spatial locations, and that matching was really trying to encode quantiles. Also it's quantized but the line chart is not.

I see a lot of potential in making the interface more actionable than descriptive. For instance, what's the effect of removing outliers or so, and comparing the cluster afterwards. Not sure if that's the purpose though

I'm not very convinced by the choice of the force layout graphs because they mutate the positions of the scatterplots or swarm for the sake of visibility. In doing so, they add some uncertainty of whether the location is accurate or not.

I liked that the outcome plots show what results are significant, but the p value could be parametrized given the different opinions of what threshold is actually significant. For the temporal trajectory, the choose of size and whether a circle is bigger than another is difficult for me to parse. I correlate the radius size with bigger being better or worse, but I believe it encodes population, and the relationship between which circle includes another circle is actually what matters. I cannot still read that well.

I liked the rule mining pcp the most, but didn't like the force layout for reasons described above. I think it could have more screen space to be able to add more rules if needed.