

1 Identifying low-dimensional trajectories of mechanically-ventilated patient systems:
2 Empirical phenotypes of joint patient+care processes to enhance temporal analysis in ARDS research

3 J.N. Stroh, Peter D. Sottile, Yanran Wang, Bradford J. Smith, Tellen D. Bennett, Marc Moss, David J.
4 Albers

5 December 14, 2023
6

7 **Abstract**

8 Mechanically ventilated patients generate waveform data that corresponds to patient interaction with
9 unnatural forcing. This breath information includes both patient and apparatus sources, imbuing data with
10 broad heterogeneity resulting from ventilator settings, patient efforts, patient-ventilator dyssynchronies, in-
11 juries, and other clinical therapies. Lung-protective ventilator settings outlined in respiratory care protocols
12 lack personalization, and the connections between clinical outcomes and injuries resulting from mechani-
13 cal ventilation remain poorly understood. Intra- and inter-patient heterogeneity and the volume of data
14 comprising lung-ventilator system (LVS) observations limit broader and longer-time analysis of such sys-
15 tems. This work presents a computational pipeline for resolving LVS systems by tracking the evolution of
16 data-conditioned model parameters and ventilator information. For individuals, the method presents LVS
17 trajectory in a manageable way through low-dimensional representation of phenotypic breath waveforms.
18 More general phenotypes across patients are also developed by aggregating patient-personalized estimates
19 with additional normalization. The effectiveness of this process is demonstrated through application to
20 multi-day observational series of 35 patients, which reveals the complexity of changes in the LVS over time.
21 Considerable variations in breath behavior independent of the ventilator are revealed, suggesting the need to
22 incorporate care factors such as patient sedation and posture in future analysis. The pipeline also identifies
23 structural similarity in pressure-volume (pV) loop characterizations at the cohort level. The design invites
24 active learning to incorporate clinical practitioner expertise into various methodological stages and algorithm
25 choices.

26 *keywords:* pulmonary ventilation; patient-ventilator asynchrony; patient-ventilator dyssynchrony; ventilator-
27 induced lung injury; respiratory distress syndrome, patient-specific modeling; knowledge representation

28 **1. Introduction**

29 Modern critical care often involves mechanical ventilation (MV) to manage patients with disorders such
30 as acute respiratory distress syndrome (ARDS), which is characterized by inflammation and pulmonary
31 edema. MV is also used to sustain highly sedate or comatose patients including those with traumatic brain
32 injury (TBI) and impaired autonomic breath control. Modern respiratory care protocols and technologies
33 [1] emphasize lung-protective strategies [2], as MV may cause ventilator-induced lung injury (VILI,[3]).
34 MV lung-protection relies on such factors as increased positive end-expiratory pressure (PEEP), decreased
35 tidal volume, or reduced driving pressure [4, 5, 6] based on understanding of lung physiology. Technological
36 advances in MV have not eliminated ventilator dyssynchrony (VD), a mismatch in patient-ventilator delivery
37 and respiratory effort timing. VD may play a role in the development and propagation of VILI, a known
38 contributor to mortality in ARDS patients [7]. Reduction in ARDS-related mortality has plateaued in
39 recent decades [8] compared with significant curtailment in the two decades prior [9]. The desire to further

40 mortality reduction motivates the continued study of MV effects as VILI and VD contribute to residual
41 ARDS mortality.

42 Clinical observables, including airway pressure (p), volume (V), and flow, are generated by the human
43 lung-ventilator system (LVS), which encompasses the dynamic interaction between patient lungs and an
44 engineered apparatus. This biomechanical LVS combines the components relevant for studying the effects of
45 mechanical ventilation (MV) on longer timescales, particularly regarding physiological derangement under
46 technological management. Non-ventilator aspects of patient care can also affect lung-ventilator dynamics
47 and associated observables. Despite the temporal richness of waveform data, analyzing them within the LVS
48 context proves challenging due to factors such as data volume (high sample rates yield millions of data points
49 per patient per day), data heterogeneity influenced by patient-specific factors and care-originating ventilator
50 setting changes, and the multi-scale nature of the problem that require consideration of intra-breath scale
51 events over extended periods to detect signatures of injuries like VILI.

52 Notable previous works [10, 11] used supervised machine learning directly on ventilator data to identify
53 the frequency and occurrence of different ventilator dyssynchronies. In addition to internal ventilator metrics,
54 the analyses also used breath properties such as peak inspiratory pressure, inspiratory-to-expiratory time
55 ratio (I:E), etc. to coarsely characterize waveform features via features familiar to practitioners [12, 13]
56 Although such descriptors facilitate operational management of respiratory care, they may be insufficient
57 to distinguish breath characteristics related to pathological lung mechanics or the timing of dyssynchronous
58 patient efforts. Identification alone does not address the evolution of breath types and effects of VD.

59 Recent approaches to LVS data analysis have focused on hybrid methods of empirical parameter fit-
60 ting [14, 15, 16] with attention to patient-ventilator dyssynchrony resolution at the waveform level. Purely
61 rule-based mechanistic models targeting specific breath features require many parameters to overcome con-
62 founding influences [16] or define models of specific VD types [14]. These research strategies have converged
63 on data-informed modeling methods as a robust tool to express waveform data through automated paramet-
64 ric representations. The present study uses a flexible model-based approach together with unsupervised ML
65 to empirically discriminate Mv breaths and begins to account for heterogeneous LVS factors. It is targeted
66 to reveal the structure and complexity of LVS evolution, and the focus on temporal factors contrasts related
67 works in ARDS research that seek to identify cohort-scale VD [10] or infer respiratory mechanics through
68 physiologic modeling [17, 14].

69 Development of relevant waveform representation models and analysis methods provide pathways for
70 informatics research to pursue minimizing VILI. Continuing toward that goal, this work presents a framework
71 for analyzing the evolution of MV breath types over extended time periods. It extends the analysis of
72 LVS behavior from the breath level to the scale of hours-to-days while considering the context of ventilator
73 settings. The method combines a model-based waveform digitization [16] with an unsupervised segmentation
74 pipeline [18], although other sufficiently flexible parametrization frameworks and variations on the theme
75 may be employed. This study's hypothesis is that respiratory behavior or other patient properties may
76 be identified from joint LVS data by separating the influence of changes in MV. Investigation proceeds by
77 examining changes in observable data that occur independently of ventilator management within the context
78 of the joint patient-ventilator system. Analysis of ARDS patient data through this perspective demonstrates
79 compact descriptions of LVS evolution, broadly categorizes MV breaths, and identifies LVS heterogeneity
80 sources that must be incorporated for further development.

81 2. Method

82 The root approach involves analyzing LVS data, including waveforms and ventilator settings, through
83 a computational pipeline that begins with model-based inference. The method projects individual LVS
84 waveform data onto personalized parametric representations and identifies patient-specific breath phenotypes
85 without consideration of sequential ordering. The evolution of LVSs may be examined through phenotypes
86 when co-labeled according to time.

87 2.1. Data

88 Mechanically-ventilated patient data were collected under the University of Colorado Multiple Institu-
89 tional Review Board (COMIRB, protocol #18-1433). These data include airway pressure, volume, and flow
90 and ventilator settings for 36 patients, all of whom had ARDS diagnoses and substantial risk of VILI as
91 featured in [19]. Children, pregnant women, and age-censored elders, and the imprisoned were excluded.
92 Esophageal pressures were recorded but used in this work; collection imposed additional exclusion criterion
93 (viz. esophageal fistula, variceal bleeding or banding, facial fracture, and recent gastric/esophageal surgery).
94 Source patients include 14 women and 22 men with median[IQR] age 59[25] years; 72% are white, 35% of
95 which identify as Hispanic or Latino. Table 1 summarizes clinical and demographic characteristics of patients.
96 Data total 1.74 million breaths over 71.14 recording-days (median 1.97[1.56] days per patient) recorded at
97 32 millisecond sampling (31.25 Hz) from Hamilton G5 ventilators (<https://www.hamilton-medical.com>).
98 Adaptive pressure and pressure-controlled ventilation modes (APVcmv and P-CVM, respectively) account
99 for 85–98% of breaths in most patients and over 94% in total. Ventilator management throughout employs
100 the ARDSnet protocols [7].

101 *Dyssynchrony labels.* An existing supervised ML technique [10, 19] identifies breath-wise VD to enrich LVS
102 evolution context and provide comparison for newly calculated labels. Type-specific VD models each label
103 breaths according to features characterizing dyssynchronous breaths (see *ibid.* SI). VD labels include normal
104 (NL), reverse triggered (RT) with early (RTe) and middle (RTm) subtypes, early flow limited (eFL) with
105 intermediate (eFLi) and severe (eFLs) subtypes, double trigger (DT) with reverse- (DTr) and patient- (DTp)
106 subtypes, and early vent termination (EVT); breath mechanics of these VDs are described in [11]. Breath
107 label vectors flag likely VD occurrence and can be summarized statistically over time intervals.

108 2.2. Windowed Model-based Inference on Individuals

109 The analysis begins with a continuous-time dynamical model that transforms observed waveform data
110 into discrete parameters via inferential methods [16]. The differential equation governing the model state y
111 is:

$$\frac{dy}{dt} + g \cdot (y(t) - y_0) = \phi(t, \omega) \quad (1)$$

112 where t is time, g is a smoothing parameter, y_0 is a reference state (such as PEEP when y represents
113 pressure), and ϕ is a time-dependent function of parameters ω . Optimizing the state y to fit observed
114 LVS waveform data over short windows yields parameters ω that encode waveform data. The relationship
115 between parameters ω and simulated state y is defined by the function ϕ . This work chooses a locally
116 periodic piecewise constant function using parameters $\omega := (a_1, \dots, a_M, \theta)$ where θ is the breath cycle length
117 determined from the data, independent of the model and M is independent hyperparameter representing

Table 1: Tabular summary of the patient cohort and associated data. ‘Monitored’ and ‘Recorded’ durations denote the number of hours spanned by data and length continuous data contents, respectively. P:F ratio is the PaO₂/FiO₂ ratio at admission, AA = African-American, AI = American-Indian, AK = Alaska, NMB = Neuromuscular Blockade

<i>Detail</i>	<i>Count</i>	<i>%</i>	<i>Median</i>	<i>IQR</i>
Monitored (hrs)			47.0	37.2
Recorded (hrs)			43.1	40
Age (years)	36		58.5	24.5
Gender				
Female	14	38.9	54.5	25.0
Male	22	61.1	58.5	26.0
Race/Ethnicity				
White	26	72.2		
Unknown/NA	5	13.9		
Black/AA	3	8.3		
AI or AK Native	1	2.8		
More than one race	1	2.8		
ARDS risk				
Pneumonia	12	33.3		
COVID	11	30.6		
Sepsis	6	16.7		
Other	3	8.3		
Pancreatitis	2	5.6		
Aspiration	2	5.6		
P:F ratio			135.9	81.0
Mortality	9	25.0		
NMB use	9	25.0		

118 the number of parameters in a . Time within the breath, defined by $\hat{t} := t - t_0 \pmod{\theta}$, is divided into M
 119 local time epochs whose lengths $\{\Delta t\}$ depend on the model resolution M , breath length θ , and partition
 120 function Υ . Each epoch is associated with an amplitude a_i , so that M determines model resolution. The
 121 piecewise-constant function ϕ can be written as

$$\phi(t; a, \theta) = \sum_{i=1}^M a_i \frac{1 - e^{-g\Delta t_i}}{g} \mathbb{I}[\Upsilon(i-1)\Delta t_i \leq \hat{t} < \Upsilon(i)\Delta t_i]. \quad (2)$$

122 The fixed function Υ apportions epoch lengths using the I:E ratio to resolve the shorter, more valuable
 123 inspiratory phase at higher resolution. Optimal parameters a are inferred from the data y^{obs} using a windowed
 124 ensemble Kalman-like smoother over short, disjoint 10-second windows of data (see [16]), although other
 125 methods suffice. The framework uses the model to infer parameters from waveform data segments and map
 126 parameters to representative waveform characterizations.

127 2.3. Pipeline

128 The computational pipeline extracts low-dimensional representations of LVS data that effectively encode
 129 relevant features of both breath waveform data and the ventilator settings associated with them. The method
 130 (Figure 1) follows [18] using model-inferred parameter distributions to uncover latent system similarities from
 131 data. The four stages of application to LVS data focus on changing system representation.

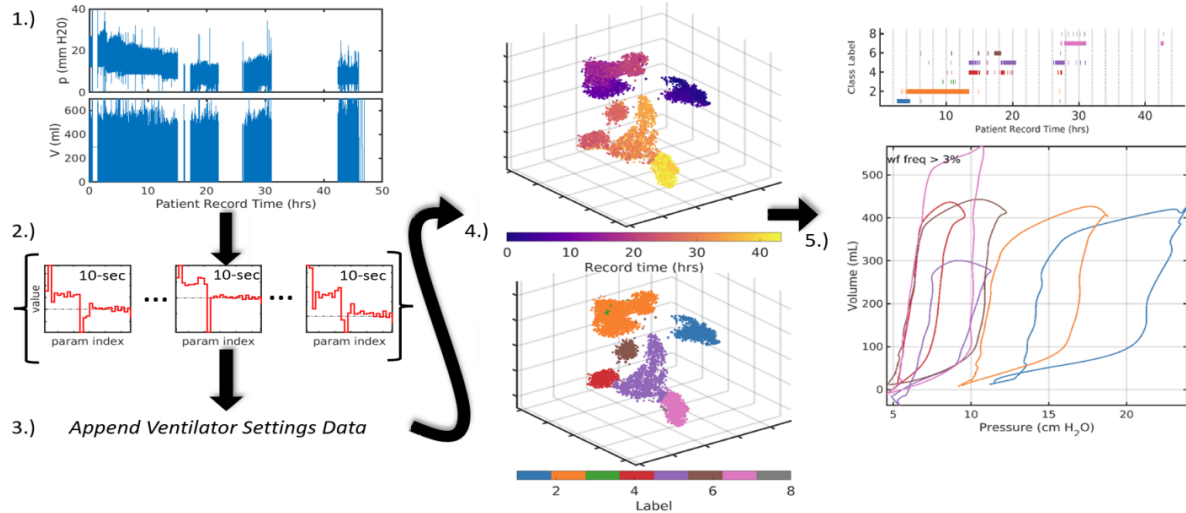


Figure 1: Broad pipeline organization. Raw data (1.) are digitally parametrized (2.) over short windows. Summaries of these vectors are computed and appended with the contextual data of ventilator settings (3.) which include information such as ventilator operation mode, PEEP, flow and pressure triggers, and minimum mandatory breath rate. Augmented LVS descriptors are projected into three dimensions using tSNE (4.) where they can be analyzed based on time ordering (top) and structural similarity via segmentation (bottom). Finally, in (5.), temporal evolution of the system is compactly encoded in the time-ordered LVS descriptor labels and their associated waveform characterizations. The process transforms raw data (1.) into a more easily comprehensible form such as (5.).

132 1.) *Waveform parametrization:*. Individual clinical records of continuous pressure (p) and volume (V) wave-
 133 forms are inferred using a model (§2.2) with moderate resolution ($M = 24$). Non-overlapping ten second
 134 windows are each encoded into M parameters by fitting the data over 1.6 second (50 points) moving sub-
 135 windows with 0.8 second overlaps (25 points, for 32 millisecond sampling). The resulting estimates are
 136 distributional samples of the 10-second windows totaling 4% fewer points than the source data.

137 2.) *Parameter Distribution Summarization:*. The parameter distributions for each interval of data aim to
 138 retain enough information to allow differentiation based on measures of relative similarity. The $2M$ param-
 139 eters are independently transformed into vectorized descriptors that collectively summarize the waveform
 140 behavior within each data window. Descriptor components include mean, quartiles, variance, and mode as
 141 well as non-gaussian measures (skewness, kurtosis, and Kolmogorov-Smirnov distance) to capture bimodal
 142 or asymmetric parameters distributions characterizing non-stationary LVS behavior. For $M = 24$, these
 143 descriptors summarize content during each 10-second interval using 38.24% less volume than the original
 144 data. The strategy reduces the temporal sampling rate (from 31.25Hz to 0.1Hz) by depicting each window of
 145 2D states as a larger vector that summarizes parameter distributions. Reduction in the overall data volume
 146 (see SIAppendix B) is governed by summary window length (under weaker stationarity assumptions) and
 147 model resolution (M).

148 3.) *Augmentation:*. Appending ventilator setting data to the parametric descriptor vectors of each window
 149 contextualizes them in the health-care process. Ventilator settings detail the mode of operation (volume
 150 control, pressure control, spontaneous, etc.), supplied targets (tidal volume) as well as various machine set-
 151 tings (trigger thresholds, ramp time, PEEP). Some ventilator settings are already represented in parametric
 152 waveform descriptors, and therefore, need not be explicitly included. For example, ARDSnet protocols bind

153 FiO₂ and PEEP ranges while realized I:E is a waveform property. Other available factors such as ventilator
154 delivery power are not considered here but may be included as needed in specific applications. Ventilator
155 mode is a nominal variable represented as set of binary variables using one-hot encoding. Interval sum-
156 maries reflect the most frequent ventilator data found for breaths in each window. Ventilator settings change
157 infrequently, and summary errors are therefore rare among estimated intervals.

158 *4.) Cluster Labeling:*. Segmentation labels groups of LVS descriptors based on content similarity and can be
159 applied at individual patient or aggregated cohort levels. Several methodological approaches can accomplish
160 this goal (e.g., see [20]). However, direct segmentation is computationally expensive when descriptors are
161 large (~ 400 -dimensional for $M = 24$). Dimensional reduction is further motivated by the desire to visually
162 inspect the quality and structure of label assignment. The t-distributed Stochastic Neighborhood Embedding
163 (tSNE,[21]) reduces high-dimensional vectors into a low-dimensional space (here, 3D) by optimizing the KL-
164 divergence between an assumed-normal distribution of the data and a t -distribution of the points in \mathbb{R}^3 . The
165 organization of embedded points approximates the local and global similarity structure [22], the targets of
166 label assignment under a given metric.

167 LVS descriptors comprise mixed variable types so the the Gower distance is a natural choice of metrics.
168 It averages over range-normalized absolute differences of continuous variables and binary dissimilarity of
169 ventilator modes (categorical variables). Uniform Manifold Approximation and Projection (UMAP,[23, 24])
170 and tSNE produce similar dimensional reductions [25] in this application (SIAppendix A.3) All individu-
171 alized results use the Matlab-native tSNE algorithm with parameters near default values (exaggeration=4,
172 perplexity=50).

173 Unsupervised learning algorithms then assigns segmentation labels to the LVS descriptors. In both
174 tSNE and UMAP LVS applications, Density-Based Spatial Clustering of Applications with Noise (DBSCAN,
175 [26, 27]) identifies groups of similar LVS descriptors from point densities in the reduced coordinate space. A
176 brief grid-search over DBSCAN parameters (min. core point neighbors 4–12; neighborhood radius 1.5–5 by
177 0.5) samples different label assignment possibilities, adopting the one that minimizes total distance between
178 cluster centroids. Experimentally, such flexible assignment sought to capture the unknown degree of variation
179 that tends to increases with the LVS record length. Use of k -means and k -medioids [28] was considered for
180 efficiency but, unlike DBSCAN, could not capture non-convex groupings that typically emerged from LVS
181 descriptors in reduced dimensions. Support vector clustering [29, 30] required too much computation time
182 to be practical for day-scale analysis.

183 *5.) Defining phenotypes.* Descriptor labels are directly associated with LVS data elements including the
184 parameter estimation windows and the waveforms contained within them. Direct interpretation of labeled
185 points is prevented by the dimensional reduction step, which embeds joint LVS descriptors into abstract,
186 similarity-determined coordinates. However, points tacitly associated with the pipeline elements used to
187 construct them, including the window times that link observations, parameters samples, summaries, and
188 tSNE coordinates The LVS data can then be analyzed based on common or central properties characterizing
189 features of each labeled group. Specifically, waveform data in a particular cluster are characterized and
190 visualized by applying the model (Eq.(2.2)) to e.g., median parameters associated with a label. These k
191 characterizations, along with their associated ventilator settings, define phenotypes of the LVS observables.

192 2.4. Phenotypes and Characterizations of LVS data

193 The phenotyping pipeline identifies elements of LVS states with similar structure, organizing short in-
194 tervals of data into discrete categories for analysis over longer timescales. Cluster labels identify LVS state
195 phenotypes of the observable data, and co-evolution of the patient-ventilator system is captured in the tem-
196 poral progression through these categories. A stated objective is to identify changes originating from the
197 patient-side of the system with no corresponding ventilator changes. These indicate the presence of con-
198 founding factors not recorded in the data such as changes in patient expectation and breathing pattern (*e.g.*,
199 patient effort, respiratory drive), lung mechanical function (*e.g.*, VILI progression or recovery from ARDS),
200 or another aspect of physiology.

201 Phenotype evolution is presented in the context of ventilator settings and in relation to VD identified
202 via [10]. Additionally, pressure-volume (pV) characterizations defined by the model image of descriptors
203 nearest to the phenotype center (*viz.* median) provide a familiar synopsis of associated waveform data for
204 each window represented in the data. Such visualizations intend to summarize key features and notable
205 changes defining the LVS trajectory.

206 Subsequent analysis and discussions employ principal component analysis (PCA), an empirical signal
207 factorization based on variance minimization [31, 32]. This tool is used to show the LVS variance occurring
208 under ventilator stationarity for qualitative analysis, as the empirically-determined basis may not represent
209 physical or relevant LVS features. Here, their intended use is to reveal the temporal structure of LVS variation
210 as these may relate to patient-side changes.

211 3. Results

212 The clinical LVS data of patients with ARDS (Table 1) is an important and practical target to test,
213 demonstrate, and document the computational phenotyping pipeline in cases which may be prone to venti-
214 lator dyssynchronies and VILI. The pipeline ties labeled breath types to specific points in time during the
215 patient record, which permits analyzing data and syntheses throughout the process. Results in this section
216 consider LVS data together with time-ordered waveform characterizations to examine LVS evolution during
217 the recorded hours of individual patients. Sequences of dyssynchrony labels generated as in [10] provide
218 additional context for exploring breath behavior.

219 Briefly, most LVS patient data are identified with 20[14] (median[IQR]) clusters using the fixed hyper-
220 parameters across the cohort. About half of these groups are infrequent and represent less than 1% of
221 the data. A median of 8[6.5] core clusters each representing more than 3% of the data account for the
222 remainder of the data of each patient. Modifying ML hyper-parameters to eliminate the low-occurrence
223 groups may consequently reduce label resolution. However, the number of labels needed to capture the main
224 LVS behaviors of each patient depend on heterogeneous factors including patient health status, the number
225 of changes in vent settings, and the total duration of the data. As recording durations span 0.7–92 hrs
226 (median[IQR] 46.8[35] hrs), the over-segmenting some LVS records to prevent loss of resolution in longer
227 ones was a preferred alternative to fully optimized individual segmentation.

228 3.1. Simple, Individual Examples

229 Figure 2 panels a–d illustrate the analysis of Patient #103 whose data consists of 7 record hours with
230 one simple ventilator setting change. Only ventilator PEEP (a) is changed while there are three primary

231 behaviors identified (b,d). The reduction of PEEP occurs about 2 hours following a rise in early flow
 232 limited breaths (eFL, panel c). This PEEP change (from 8 to 5 cmH₂O) shifts peak pressure from 16 to 12
 233 cm H₂O for about an hour, at which time higher esophageal pressures returns. These breaths are identified as
 234 normal (NL) [10]. Increased specificity may be pursued by local segmentation or other dimensional reduction
 235 methods.

236 *A closer look at label 1 of patient #103:*. The first principal component loadings (panel e, black) for LVS
 237 descriptors over the first 5-hour period track the sequence of normal and eFL VD labels (f, shown as 5
 238 minute statistics for clarity). Within the same breath phenotype (label 1), the *sign* of the component
 239 loading statistically the eFL VD labels (AUROC=0.8718); high positive values are associated with eFL
 240 breaths (f,g; green) where pressure maxima proceed volume maxima. These LVS variations result from
 241 changes in the patient component, as there is no change of ventilator settings. Note that direct correlation
 242 between continuous loading values on 10 second windows and statistical breath-wise binary VD label is not
 243 well-defined while binary-to-binary comparison is.

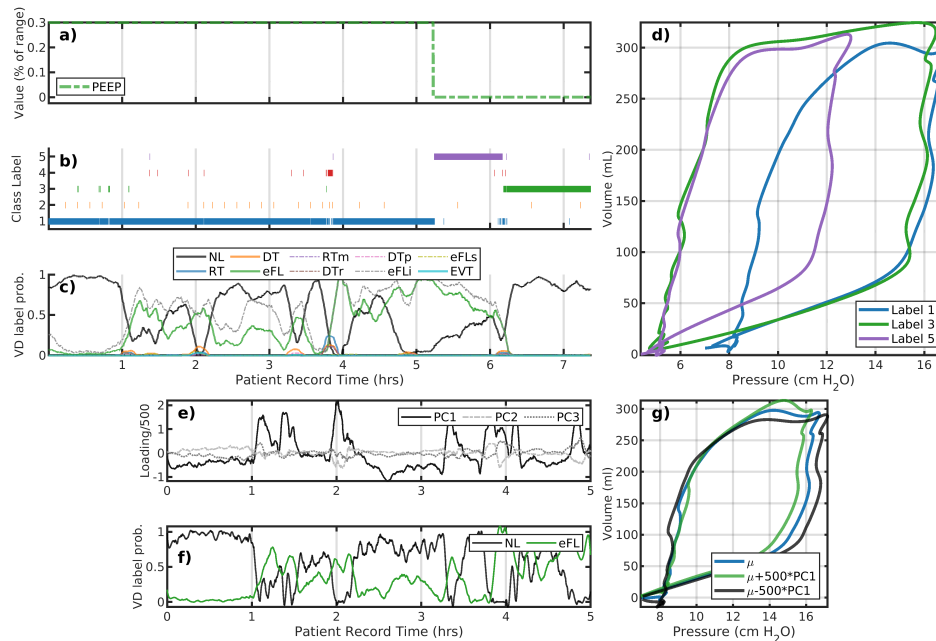


Figure 2: Analysis of patient #103 LVS data (a–d) and the initial 5-hour interval (e–g). Panels a–c correspond to changes in ventilator settings, segmentation labels, and identified VD type, respectively. The horizontal axis for these panels is the patient record time in hours. The panel (d) shows the model image of segmented data median parameters, which characterize the pV loops of breaths with that label (shown with the same color). Evolution of the LVS can be parsed pictorially from these figures. Large positive variations in the first principal component loading (e, black) for the initial 5-hour period align with VD labels indicating eFL type breaths (f) for this period. Specifically, this suggests discrimination of breaths shapes (g) can be differentiated using qualitatively criterion on local loadings or other segmentation.

244 The patient #113 (Figure 3) dataset is nearly twice as long with again only one PEEP change occurring
 245 after 10.5 hours of the 15.6 hour record. Breathes are stably identified as normal-type until about 8 hours,
 246 occupying two cluster-identified similar breath shapes. This is followed briefly by eFL breaths and a transition
 247 to a new characterization (label 8, light green) for about 30 minutes. In the following period (9–14 hours),
 248 breaths are characterized by lower pressure maxima (label 10, gold); these are associated/identified with

249 reverse-trigger breaths (primarily RTm) and waveforms featuring pronounced inspiratory pressure drop.
250 The reduction in PEEP slightly increases the incidence of normal breaths during 11–14 hours although this
251 results in the more frequent appearance of shallow breaths (label 13, red).

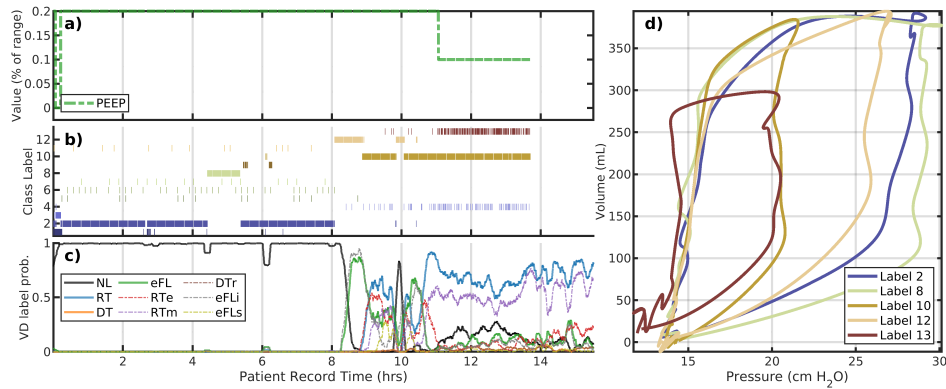


Figure 3: The patient #113 evolution also includes only PEEP changed. The layout is the same as panels a–d of the previous figure. Under constant ventilator settings, breaths undergo transition several times including intervals of VD prior to PEEP change around 10.5 hours. A 1-hour long shift from label 2 to 8 occurs around 8 hours during which breaths decrease peak pressure and includes an increase in eFL and RT VD occurrence. After the PEEP change, breaths remain highly dyssynchronous and primarily centered around the characterization with label 10.

252 3.2. More complex LVS evolution examples

253 Cases presented in the previous subsection are atypical in that patient records in the data set are typi-
254 cally longer (>24 hours), include many ventilator settings changes, and segment into a larger collection of
255 phenotypic breaths. Figure 4 illustrates the analysis of patient #114 whose LVS undergoes multiple changes
256 over a 24-hour data period. The portion during 7–14.5 hours is dominated by normal breaths that spans
257 two labels with similar characterizations as pV loops but differ in mean respiratory rate. The difference is
258 minor (the mean difference is less than 20 milliseconds), although this affects model parameter and could
259 combined via posterior analysis, with small DBSCAN hyper-parameter changes, or coarser period binning.
260 Changes in flow trigger settings occur around 3 hours and reduce the occurrence of eFL near the star of the
261 record, associated with caving in pV loops (label 1, dark blue). Dyssynchronies return when the flow trigger
262 is returned to its initial value, near 15 hours. PEEP and tidal volume targets are also adjusted several times.
263 Brief ventilator changes in ventilator mode around 20 and 23 hours allow spontaneous breathing which have
264 a profoundly different pV characterizations (label 12, tan). The interim period (20.5–22.5 hours) consists of
265 primarily normal breaths (label 13, brown) under the default pressure-control mode.

266 *A closer look at Label 10 of patient #114.* Breath phenotype analysis of patient #114 indicates no ventilator
267 setting changes during the record interval 15–21 hours. Although one phenotypic breath dominates this pe-
268 riod, ML-labeled dyssynchronies intimates much more variability. Principal components during this interval
269 (Fig.4b) suggest that the evolution is irregular. While pressure characterizations suggest the differences are
270 largely attributed to pressure and inspiration duration, full characterizations indicate that breaths in this
271 period are very heterogeneous in pV relationships. The *continuous* evolution through these subtypes – and
272 their comparative differences to the other types – leads to their identification as a consistent group. SI Figure
273 A.7 provides another case using principal components to further differentiate breath types with implications.

It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).

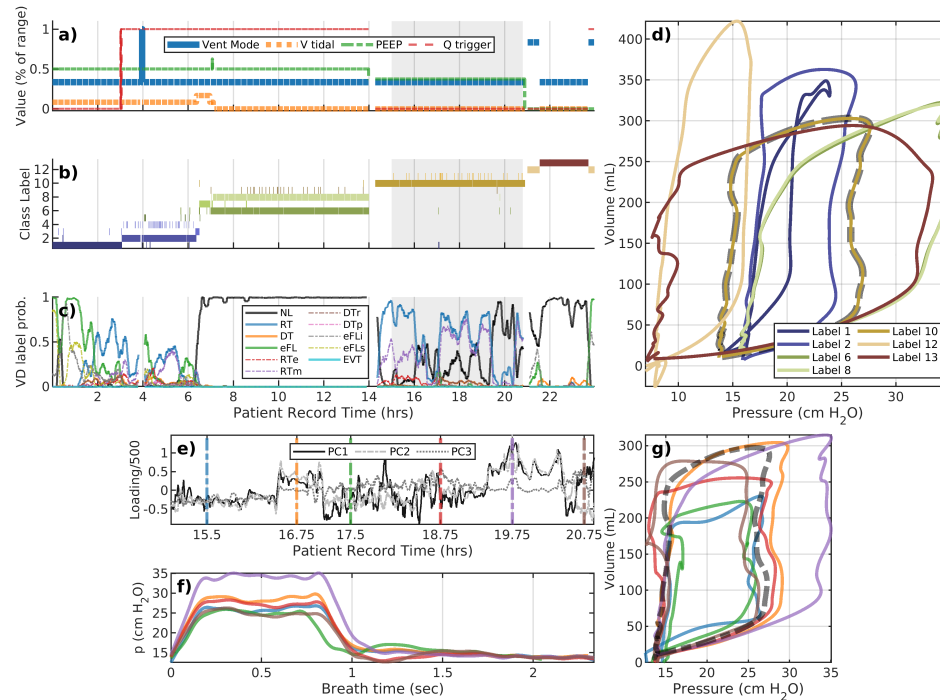


Figure 4: A representative example: patient #114. The upper plot layout is the same as the previous figure. The lower plot examines the variability, and a shortcoming of low resolution segmentation to capture changes that may be highly diverse at a local level. The mean – the dashed black line – coincides with the golden pV loop (cluster #10) in the upper plot. The many distinct breath subtypes identified are more similar than to other main types in the upper plot; as a result, they are grouped together at this choice of hyperparameters.

274 Figure 5 features data of the patient #149 LVS has little identified dyssynchrony. While most breaths as
 275 identified as normal, the system evolution diagram indicates irregularities in breath properties. In particular,
 276 label 1 (dark blue) regularly present under multiple PEEP settings in the pressure-controlled volume targeting
 277 mode, and are intermixed with other labels (e.g., 3,5, and 7) whose waveform characterizations are dissimilar.
 278 System heterogeneity within the patient makes parsing the evolution more complicated, as spontaneous
 279 breathing is possible during 12–16 hours and 23–24 hours under different PEEP values. Nevertheless, the
 280 space of breaths is considerably reduced in labels.

281 3.3. Cohort scale phenotypes of breath shapes

282 A cohort level breath characterization is achieved by further segmenting the population of individuals
 283 characterizations in non-dimensional form. The phenotyping pipeline yields a total of 721 patient-specific LVS
 284 characterizations across the cohort. Attempts to directly cluster these full characterizations was unsuccessful.
 285 Secondary tSNE-DBSCAN grouping could not identify cross-patient commonalities of these characterizations
 286 due to LVS-specific heterogeneities such as tidal volume (a target set by patient sex, height), PEEP (e.g.,
 287 whether even or odd values are used), and respiratory rate (patient and sedation dependent). However, these
 288 factors may be accounted for by segmentation of pV loop shape rather than full LVS behavior. Pressure
 289 and volume characterizations are sampled in pV-space and then translated and scaled into the range [0,1].
 290 The pV normalization accounts for differences in PEEP, tidal volume, and respiratory rate while preserving
 291 the differences in the pressure-volume relationship to extract the phenotypic shapes of breaths occurring

It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).

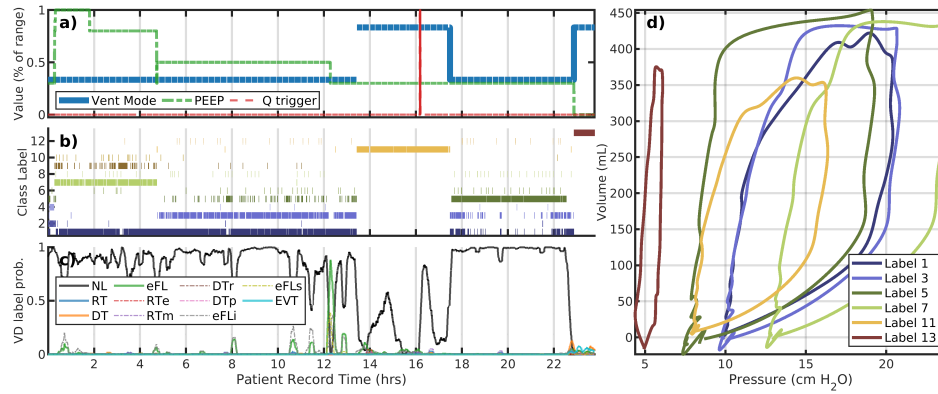


Figure 5: The LVS evolution of patient 149 label#1 is discontinuous in time and occurs under different PEEP values. There is a lot of waveform variation present within the largely VD-less evolution, and significant changes in non-ventilator aspects of the LVS.

292 throughout the estimated entire dataset. Using this method, a cohort-scale tSNE-DBSCAN analysis identifies
293 20 breath shape phenotypes along with a collection of 27 outliers. Figure 6 illustrates this normalized re-
294 organization along with the median pV representatives and pressure traces for each identified group.

295 Meta-characterization depends on several hyper-parameters associated with tSNE and DBSCAN which
296 influence label granularity as well as thresholds defining outlier groups (SI Fig. A.8). Selected parameters
297 aimed to maximize the number of phenotypes while minimizing the number of outliers with number of labels
298 easily presented in an array; the results are qualitatively similar for nearby parameters. Table 2 summarizes
299 the occurrence and contents of this grouping.

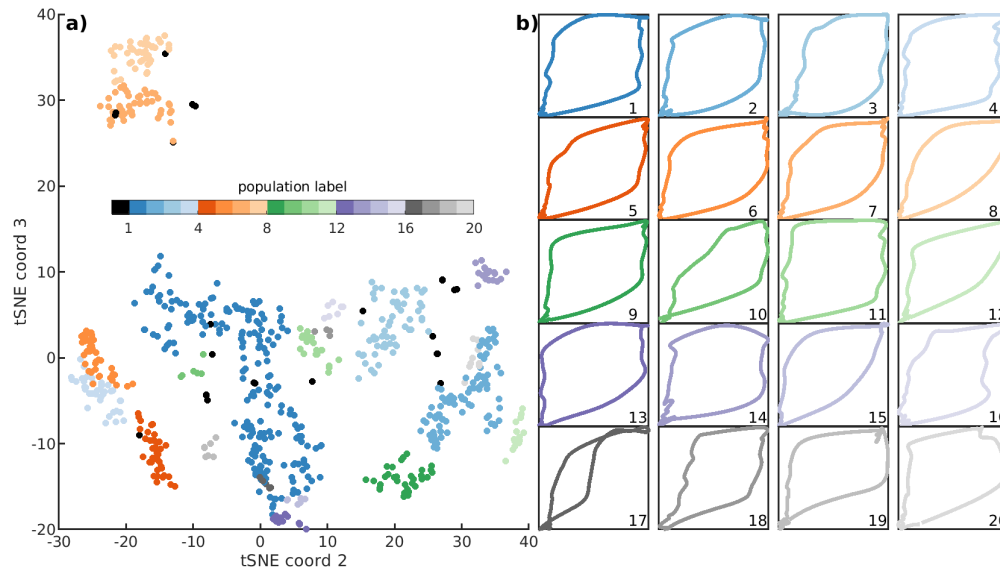


Figure 6: Cohort scale breath types via segmentation of batched individual data. In panel a, each dots corresponds to one of $N = 721$ normalized pV loop characterizations extracted at an individual level. Labels (colors) identify 20 fundamental breath shapes comprising the dataset of 1.74M breaths. Un-grouped individual characterizations correspond to 27 outliers (black). The 20 pV trace shapes nearest to groupwise medians are shown at right (b) (see SI Fig. A.8 for outliers). Note that the xy -axes are the normalized pressure and volume in non-dimensional form.

Table 2: Tabulated occurrence of labels within the cohort. Columns indicate the cohort label (with 0 indicating the group of unrelated outliers), occurrence percentage, the number of patients represented in the label, and the number of individual LVS phenotypes comprising each label. Eleven non-outlier groups with occurrences greater than $>2.5\%$ constitute 87% of the data with outliers (Label 0) occupying an additional $\sim 3.5\%$. Two breath types (*italics*) contain only phenotypes of patient #141, whose average inspiration compliance is consistently ~ 7.5 ml/cm H₂O) and distinct from the 1.5–3 ml/cm H₂O range typical across other patients.

Label	Percent	N_{pat}	N_{pheno}	Label	Percent	N_{pat}	N_{pheno}
0	3.48	14	27	-	-	-	-
1	29.68	28	200	11	2.51	6	20
2	11.10	12	85	12	2.19	1	16
3	9.74	13	58	13	1.69	2	14
4	6.52	4	38	14	1.69	5	18
5	5.51	4	39	15	1.43	1	6
6	5.31	7	39	16	1.29	2	7
7	5.30	10	48	17	0.88	4	5
8	4.74	3	43	18	0.19	4	4
9	3.66	4	32	19	0.10	3	7
10	2.93	4	6	20	0.07	4	9

300 3.4. Synthesis

301 Variations in pressure and volume observations of MV patients result from ventilator setting changes and
302 patient dynamics. To understand the evolution of this system, one must jointly consider both patient and care
303 processes. Analyzing patient state through breath data, especially for VILI detection and to track ARDS
304 progression, requires considering ventilator settings. LVS evolution is primarily influenced by ventilator
305 setting changes (e.g., PEEP, mode, tidal volume), with secondary changes indicative of patient progression
306 or non-ventilator care. Analyzing periods of ventilator stationarity showed local breath evolution unrelated
307 to MV changes, and empirically-identified breath variations agreed with ML-derived statistical labels of VD
308 (patient 103 in Fig 2; patient 111 in Fig SI A.7). Local analysis or sub-phenotyping may resolve cases with
309 more complicated evolution (e.g., patient 114 in Fig. 4b, although additional factors such as sedation and
310 patient restfulness may be required to characterize and differentiate them.

311 Cohort-level segmentation of LVS behavior phenotypes was confounded by heterogeneities in patients as
312 well as ventilator settings that depend on patient-specific properties. Specifically, tidal volume and breath
313 rate differences among and within patients could not be accounted for, leading to partitioning of both patients
314 and LVS behavior. However, secondary clustering of normalized 721 individual pV characterizations yielded
315 20 pV shapes and a collection of 27 outliers. As ventilator settings and breath rate factor into individual
316 LVS segmentation, resulting pV characterization phenotypes inherit some aspect of those data indirectly.

317 Results and following conclusions naturally depend on algorithm hyperparameters that should be chosen
318 in relation to targeted applications. The proposed methodology supports clinical expert guidance in selecting
319 guidance of application, identifying specific features to be resolved in models, and phenotype interpretation;
320 each stage of Fig. 1 supports expert-in-the loop refinement. Use of the phenotyping pipeline and subsequent
321 analysis leads to methodological conclusions regarding application to real, clinical LVS data:

- 322 1. In application to a cohort of ARDS patient clinical data, system evolution is much easier to visualize,
323 track, assess, and understand when the LVS is represented in a discrete, low-dimensional form as an
324 evolution of phenotypes and waveform characterizations.

- 325 2. LVS-individual phenotypes depend on hyper-parameters of the dimensional reduction and labeling
326 stages. Parameters affect the granularity of resolved phenotypes, and should be chosen flexibly to
327 account for the desired quality of LVS temporal resolution as well as the length and complexity of the
328 target dataset.
- 329 3. Sub-segmentation of phenotypes is possible, and may be used to align phenotype definitions with
330 additional information streams such as VD identification. Cohort scale analysis is also possible from
331 batched individual phenotyping to provide a coarser but unified basis for analyzing common trends
332 and evolution of LVSs.

333 4. Discussion

334 Research into ARDS and VILI involves studying patient-ventilator interactions and may benefit from
335 representation of the lung-ventilator system (LVS) over time. Typically, the available data include airway
336 pressure, volume (or flow), and ventilator settings. This study introduces a framework to transform these
337 LVS data into meaningful, low-dimensional characterizations of LVS state that facilitate analysis of LVS
338 behavior and its evolution during patient care. The process involves aggregating segmented analyses of
339 individual patient over short (10-second) intervals. Consequently, the observable LVS data is condensed into
340 a small set of patient-level phenotypes, making it discrete and more manageable compared to continuous
341 high-resolution waveforms and breath-wise ventilator settings.

342 Experiments conducted on clinical data of 35 patients with strong ARDS risks, including 8 2020 pa-
343 tients with COVID-19, found the automated phenotyping process sufficient to discern between changes in
344 the ventilator and the patient components of the LVS system. Individual LVS phenotypes were primarily
345 determined by ventilator setting changes, given that changes in mode, PEEP, and tidal volume profoundly
346 affect waveform shapes. However, temporal changes in phenotype uncoordinated with ventilator changes
347 were also present in nearly all patients with more than 12 hours of data, revealing changes in the patient
348 side of the LVS system. Unlike the rapid and instantaneous transitions related to MV changes, patient-
349 side changes were often a more continuous but non-monotonic progression together with transient behavior.
350 These behaviors could be detected through principal component analysis of data over intervals of static MV,
351 but additional EHR data required to adequately explain them are presently unavailable.

352 In cases with limited complexity, LVS phenotypes corresponded well with ML-labeled VD, as appear re-
353 lated to ventilator setting and behavior in response to changes. For more complex cases, breath phenotypes
354 did not consistently differentiate normal and dyssynchronous breaths (Figs. 3,4). Signatures of VD can be
355 subtle but remain discernible through empirical analysis (such as PCA) or sub-segmentation of individual
356 phenotypes, and in the original data associated with each phenotype. Nevertheless, an important consid-
357 eration for future work is the optimization of phenotype resolution, which is defined by of easily tunable
358 hyper-parameters (viz. those of tSNE, UMAP, and DBSCAN) as well as deeper factors discussed below.

359 While granular refinement is accessible through targeted sub-segmentation of LVS phenotypes, batch
360 analysis makes convenient a cohort analysis to assess the frequency of different breath shapes. For specific
361 choice of hyper-parameters, the ~ 1.5 M breaths reduce to a small set of 20 pV loop shapes and a set of
362 27 outliers. Signs of dyssynchrony are apparent in these core shapes such as ineffective triggering (10,18),
363 and mild and severe flow limitation or patient effort (14,16, and 17 respectively). (Also, types 19 and 20
364 qualitatively suggest RT with late insufflation pressure drops indicative of patient effort.) This qualitative

365 VD classification remains formally unvalidated in this work, as cohort scale breath shapes mask important
366 considerations such as PEEP, tidal volume, and respiratory rate. Additionally, esophageal pressures were not
367 encoded into phenotypes but are required to identify certain VD types [11]. Although individual phenotypes
368 were often too coarse to differentiate normal and dyssynchronous breaths, this coarser segmentation could
369 be refined by scaling to larger cohorts and longer patient data series.

370 *4.1. Method Choices, Variations, and phenotype consolidation*

371 The pipeline presented makes several choices regarding application-specific details. Algorithm choices
372 such as model resolution, estimation/summary window length, the intra-breath partition (viz. Υ of Eqn.(2)),
373 and omission of certain data are needed to balance efficiency with the quality of results. These choices
374 were influenced by practical consideration such as stationarity over 3–4 breaths and model consistency [16].
375 Many modifications and changes to latter stages of the labeling process are possible including: hierarchical
376 analysis of specific times at individual or cohort scales, characterizing breath types occurring under particular
377 ventilator settings or modes, or incorporating other factors such as patient sedation level.

378 It remains essential for certain applications to examine breath shape independently of ventilator settings.
379 For example, investigating signatures of ventilator dyssynchrony [11] may require normalizing breath features
380 to account for ventilator settings that affect pressure and volume waveform maxima. To account for PEEP
381 and tidal volume in this process, normalization during data segmentation requires a similarity measure to
382 be invariant under translation and scaling, respectively. Use of parametrized waveform descriptors does not
383 eliminate this problem. Circumventing these obstacles is possible by comparing waveform characterizations
384 and merging labels based on characterization similarity, gauged by the difference between normalized char-
385 acterized pV loops. This approach, explored in §3.3, applies to experiments that differ significantly in scale,
386 and is motivated by the desire to link uncontrolled human LVS data with in vivo animal experiments (e.g.
387 [33]).

388 *4.2. Limitations and Improvements*

389 Combining data assimilation-based parametrization with unsupervised learning ([18]) overcomes primary
390 shortcomings of existing approaches. In particular, the mechanism-free encoding of waveform data into
391 parameters with a priori definition circumvents patient- and care-dependent heterogeneity which strongly
392 limit physiological model use in this domain ([16]). This greatly alters the representation LVS system: the
393 rapid temporal sampling in two dimensions (p,V) is transformed into a low-frequency sampling of model
394 parameter distributions (SIAppendix B) under stationarity assumptions. Accounting for irregularity in
395 sample dimensions (viz. the number of points representing each breath) caused by variable respiratory rate
396 in these breath-wise analyses is unnecessary in the continuous-time windowed approach.

397 An important limitation of this work and its clinical application regards the dependence on hyper-
398 parameters and a distance function used in dimensional reduction and group labeling. Fixed tSNE param-
399 eters and a very narrow DBSCAN parameter range not adequately account for individual record length,
400 internal waveform heterogeneity, or the number of unique ventilator settings. Examples showed that chosen
401 parameters of the segmentation processes were insufficient to produce phenotypes that corresponded with
402 VD types using a known method, while also generating many smaller, low occurrence phenotypes. Selecting
403 algorithm parameters to align identified phenotypes with VD labels is likely achievable as an application-
404 specific optimization. This topic is important for clinical use but lies beyond the present scope focused on

405 low-dimensional representation of LVS evolution. Additionally, the uniform weighting of components in the
406 similarity metric used for dimensional reduction may not be optimal. Strategic weighting would require
407 objective criteria compatible with mixed variables to apportion LVS descriptors correctly. However, a well-
408 specified metric may improve low-dimensional tSNE or UMAP representation so that estimating VD severity
409 is feasible within the LVS phenotyping results.

410 Algorithmically-defined LVS phenotypes did not include several important factors that limit the strength
411 of conclusions about them. The pipeline process ignored esophageal pressure data because their rarity
412 limits generalizability, their highly-localized features require high resolution to capture parametrically, and
413 record inconsistencies (gaps, calibration) prevent continuous time characterization. Exclusion of this variable,
414 essential to defining certain types of VD [11], limits phenotype ability to distinguish certain types of VD
415 from airway observations. In addition, the model parameter definitions relies on ventilator-identified breath
416 rate, and the pipeline therefore lacks the flexibility needed to identify double-triggered VD events that occur
417 over multiple ventilator cycles.

418 Most importantly, the analysis did not consider extra-LVS influences on observable data, such as patient
419 sedation, neuromuscular blockade use, posture, and airway moisture and secretions. These patient-state
420 variations undoubtedly impact observed data and must be included to properly vet phenotypes identified
421 under ventilator stationarity (cf. Figs. 2e–g and 4e–g).

422 *Appropriate Normalization.* Cohort-scale segmentation of Sec3.3 normalized pressure and volume to a stan-
423 dard interval for inter-comparison. This waveform rescaling depends on local tidal volume, peep, and driv-
424 ing pressure which could be included as feature components to improve discrimination. Scaling volumes
425 by predicted body weight accounts for patient invariants (sex and height) while pressure scaling remains
426 co-dependent patient+care processes (viz. plateau pressure and assigned PEEP, assigned tidal volume per
427 kg in adaptive pressure modes). Waveform data and their parametric summaries dominate the dimension
428 of feature vectors. Unreported experiments using naive normalization yielded cohort breath labels that seg-
429 mented patients rather than refining breath types groupings. The topic warrants investigation to identify
430 appropriate cohort scale normalizations in addition to metrics and weights needed to balance the roles of
431 normalized waveforms and associated scaling factors in label identification.

432 *4.3. Concluding Remarks*

433 This work demonstrates an effective operationalization of lung-ventilator systems for analyzing patient-
434 ventilator interactions and breath types over extended timescales to facilitate the study of VILI and its
435 connection VD. Computationally defined phenotypes consolidate LVS states into classes, reducing patient-
436 ventilator dynamics to evolution of discrete phenotypic states. The development permits investigation of
437 time-dependent changes in MV patients within the context of applied care from observable data. The
438 approach encourages hypothesis formulation regarding the role of VD and MV duration on VILI by preserving
439 time-ordered links between LVS data and low-dimensional representations that are easier to analyze and
440 study. The pipeline organization is structured around active learning to incorporate domain expert knowledge
441 into waveform feature targets, ventilator setting inclusion, and group similarity definitions.

442 The hybrid method incorporates model-based data assimilation and unsupervised machine learning to
443 simplify LVS data into empirically-grouped rule-based descriptors. A suitable next step for understanding
444 LVS evolution is the use of symbolic dynamics [34, 35, 36] to examine and identify common temporal patterns

445 arising within patient cohorts. An initial step toward this goal is systematizing individual LVS evolutions
446 within cohort-scale phenotypes (cf. §3.3) and tuning hyper-parameters for the resolution needed target this
447 specific research goal.

448 **Acknowledgements**

449 This work is supported by National Heart Lung and Blood Institute awards 5R01HL151630 “Predict-
450 ing and Preventing Ventilator-Induced Lung Injury” (JNS, BJS, DJA) and K23HL145011 “The Detection,
451 Quantification, and Management of Ventilator Dyssynchrony” (PDS). Big-ups as always to Meg Rebull for
452 local administrative support.

453 **Declarations of Interest**

454 None. The authors have no conflicts of interest to disclose.

455 **References**

- 456 [1] Eddy Fan, Jesus Villar, and Arthur S Slutsky. Novel approaches to minimize ventilator-induced lung
457 injury. *BMC medicine*, 11(1):1–9, 2013.
- 458 [2] Gerard F Curley, John G Laffey, Haibo Zhang, and Arthur S Slutsky. Biotrauma and ventilator-induced
459 lung injury: clinical implications. *Chest*, 150(5):1109–1117, 2016.
- 460 [3] Arthur S Slutsky and V Marco Ranieri. Ventilator-induced lung injury. *New England Journal of*
461 *Medicine*, 369(22):2126–2136, 2013.
- 462 [4] Roy G Brower and Gordon D Rubenfeld. Lung-protective ventilation strategies in acute lung injury.
463 *Critical care medicine*, 31(4):S312–S316, 2003.
- 464 [5] Nicola Petrucci and Walter Iacovelli. Lung protective ventilation strategy for the acute respiratory
465 distress syndrome. *Cochrane Database of Systematic Reviews*, (3), 2007.
- 466 [6] Yuda Sutherland, Maria Vargas, and Paolo Pelosi. Protective mechanical ventilation in the non-injured
467 lung: review and meta-analysis. *Annual Update in Intensive Care and Emergency Medicine 2014*, pages
468 173–192, 2014.
- 469 [7] Acute Respiratory Distress Syndrome Network. Ventilation with lower tidal volumes as compared with
470 traditional tidal volumes for acute lung injury and the acute respiratory distress syndrome. *New England*
471 *Journal of Medicine*, 342(18):1301–1308, 2000.
- 472 [8] Shea E Cochi, Jordan A Kempker, Srinadh Annangi, Michael R Kramer, and Greg S Martin. Mortality
473 trends of acute respiratory distress syndrome in the united states from 1999 to 2013. *Annals of the*
474 *American Thoracic Society*, 13(10):1742–1751, 2016.
- 475 [9] Marc Moss and David M Mannino. Race and gender differences in acute respiratory distress syndrome
476 deaths in the united states: an analysis of multiple-cause mortality data (1979–1996). *Critical care*
477 *medicine*, 30(8):1679–1685, 2002.
- 478 [10] Peter D Sottile, David Albers, Carrie Higgins, Jeffery Mckeehan, and Marc M Moss. The association
479 between ventilator dyssynchrony, delivered tidal volume, and sedation using a novel automated ventilator
480 dyssynchrony detection algorithm. *Critical Care Medicine*, 46(2):e151, 2018.

- 481 [11] Peter D Sottile, David Albers, Bradford J Smith, Marc M Moss, et al. Ventilator dyssynchrony–
482 detection, pathophysiology, and clinical relevance: A narrative review. *Annals of Thoracic Medicine*,
483 15(4):190, 2020.
- 484 [12] G Schmidt. Ventilator waveforms: clinical interpretation. *Principles of Critical Care*. New York:
485 *McGrawHill*, 427:443, 2005.
- 486 [13] Elizabeth Emrath. The basics of ventilator waveforms. *Current pediatrics reports*, 9:11–19, 2021.
- 487 [14] Deepak K Agrawal, Bradford J Smith, Peter D Sottile, and David J Albers. A damaged-informed lung
488 ventilator model for ventilator waveforms. *Frontiers in physiology*, 12, 2021.
- 489 [15] Cong Zhou, J Geoffrey Chase, Qianhui Sun, Jennifer Knopp, Merryn H Tawhai, Thomas Desaive, Knut
490 Möller, Geoffrey M Shaw, Yeong Shiong Chiew, and Balazs Benyo. Reconstructing asynchrony for
491 mechanical ventilation using a hysteresis loop virtual patient model. *BioMedical Engineering OnLine*,
492 21(1):1–20, 2022.
- 493 [16] JN Stroh, Bradford J Smith, Peter D Sottile, George Hripcsak, and David J Albers. Hypothesis-driven
494 modeling of the human lung-ventilator system: A characterization tool for acute respiratory distress
495 syndrome research. *Journal of Biomedical Informatics*, page 104275, 2022.
- 496 [17] Michelle M Mellenthin, Siyeon A Seong, Gregory S Roy, Elizabeth Bartolák-Suki, Katharine L Ham-
497 ington, Jason HT Bates, and Bradford J Smith. Using injury cost functions from a predictive single-
498 compartment model to assess the severity of mechanical ventilator-induced lung injuries. *Journal of*
499 *Applied Physiology*, 127(1):58–70, 2019.
- 500 [18] Y Wang, JN Stroh, George Hripcsak, Cecilia C Low Wang, Tellen D Bennett, Julia Wrobel, Caroline
501 DerNigoghossian, Scott Mueller, Jan Claassen, and DJ Albers. A methodology of phenotyping icu
502 patients from ehr data: high-fidelity, personalized, and interpretable phenotypes estimation. *in review*
503 *Journal of Biomedical Informatics*, xx(x):xxx, *in review* 2023.
- 504 [19] Peter D Sottile, Bradford Smith, Marc Moss, and David J Albers. The development, optimization, and
505 validation of four different machine learning algorithms to identify ventilator dyssynchrony. *medRxiv*,
506 2023.
- 507 [20] Mahamed GH Omran, Andries P Engelbrecht, and Ayed Salman. An overview of clustering methods.
508 *Intelligent Data Analysis*, 11(6):583–605, 2007.
- 509 [21] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of machine*
510 *learning research*, 9(11), 2008.
- 511 [22] George C Linderman and Stefan Steinerberger. Clustering with t-sne, provably. *SIAM Journal on*
512 *Mathematics of Data Science*, 1(2):313–332, 2019.
- 513 [23] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projec-
514 tion for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- 515 [24] Connor Meehan, Stephen Meehan, and Wayne Moore. Uniform manifold approximation and projection
516 (UMAP v4.2). *MATLAB Central File Exchange*, 2022. [https://www.mathworks.com/matlabcentral/
517 fileexchange/71902](https://www.mathworks.com/matlabcentral/fileexchange/71902).
- 518 [25] Dmitry Kobak and George C Linderman. Initialization is critical for preserving global data structure
519 in both t-sne and umap. *Nature biotechnology*, 39(2):156–157, 2021.
- 520 [26] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for
521 discovering clusters in large spatial databases with noise. In *KDD*, volume 96, pages 226–231, 1996.

- 522 [27] Erich Schubert, Jörg Sander, Martin Ester, Hans Peter Kriegel, and Xiaowei Xu. DBSCAN revisited,
523 revisited: why and how you should (still) use DBSCAN. *ACM Transactions on Database Systems*
524 (*TODS*), 42(3):1–21, 2017.
- 525 [28] Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements of*
526 *statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
- 527 [29] Asa Ben-Hur, David Horn, Hava T Siegelmann, and Vladimir Vapnik. Support vector clustering. *Journal*
528 *of machine learning research*, 2(Dec):125–137, 2001.
- 529 [30] Jaewook Lee and Daewon Lee. Dynamic characterization of cluster structures for robust and in-
530 ductive support vector clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
531 28(11):1869–1874, 2006.
- 532 [31] Harold Hotelling. Analysis of a complex of statistical variables into principal components. *Journal of*
533 *educational psychology*, 24(6):417, 1933.
- 534 [32] C Radhakrishna Rao. The use and interpretation of principal component analysis in applied research.
535 *Sankhyā: The Indian Journal of Statistics, Series A*, pages 329–358, 1964.
- 536 [33] H Farooqi, AM Sosa, PD Sottile, DJ Albers, and BJ Smith. Experimentally induced ventilator dyssyn-
537 chrony increases injury during prolonged ventilation of endotoxin-injured mice. In *C72. HOUSE OF*
538 *ARDS... AND MECHANICAL VENTILATORY SUPPORT*, pages A5780–A5780. American Thoracic
539 Society, 2023.
- 540 [34] José M Amigó, Karsten Keller, and Valentina A Unakafova. Ordinal symbolic analysis and its application
541 to biomedical recordings. *Philosophical Transactions of the Royal Society A: Mathematical, Physical*
542 *and Engineering Sciences*, 373(2034):20140091, 2015.
- 543 [35] Douglas Lind and Brian Marcus. *An introduction to symbolic dynamics and coding*. 2nd edition, 2021.
- 544 [36] Yoshito Hirata and José M Amigó. A review of symbolic dynamics and symbolic reconstruction of
545 dynamical systems. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 33(5), 2023.

546 **Appendix A. Supporting Figures**

547 *Appendix A.1. Intracluster normal and eFL in p111, label2*

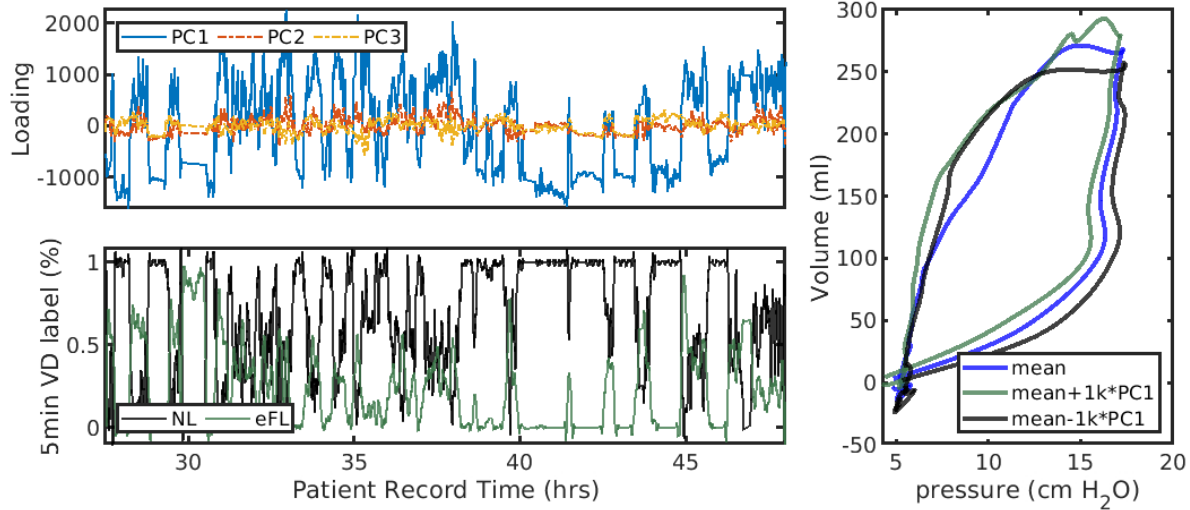


Figure A.7: The sign of PC1 loading roughly divides the VD classes in p111, label2. A threshold for the PC1 loading at zero roughly separates NL and eFL labels by 34%/65% and 85%/14%, respectively, with NL labels strongly associated with negative loadings. The optimal threshold (~ 0.05) offers only subtle improvement. The right panel illustrates low fidelity changes in the cluster median pV loop (blue) when modified by these negative (black, more associated with NL) and positive (green, eFL) loadings. Note that this involves comprising 10-second properties (representing typically ~ 3 breaths) to breathwise labels, and some representation errors thus arise from summarizing binary VD labels distributionally over all breaths intersecting a 10-second analysis window.

548 *Appendix A.2. Outlier individual cluster characterizations in the cohort segmentation*

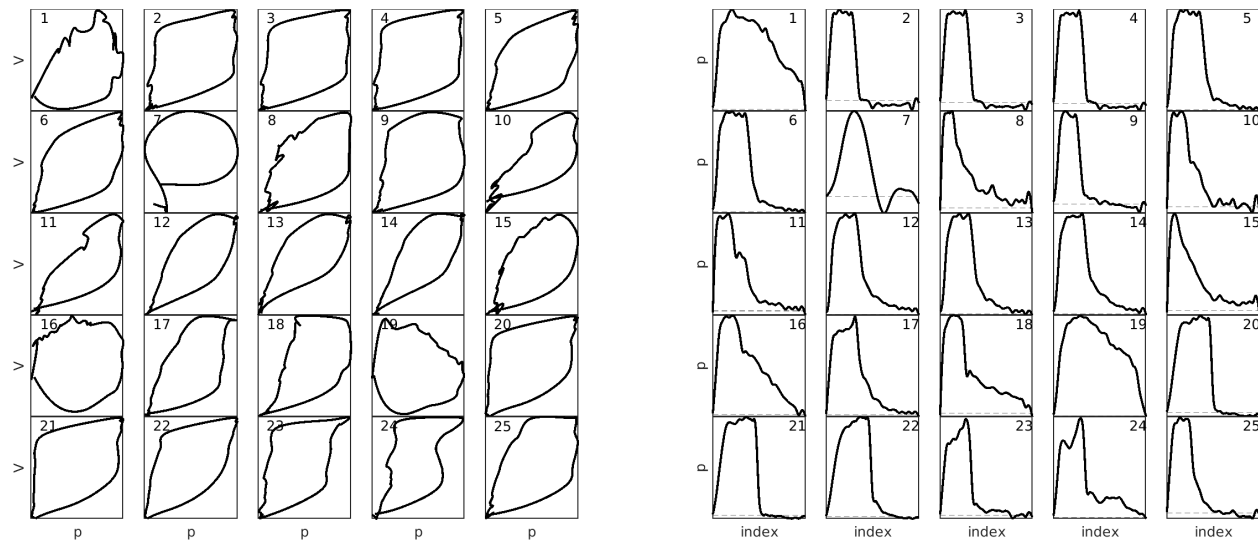


Figure A.8: Outlier pV (left) and pressure waveform (right) characterizations from two-stage cohort LVS phenotyping, shown in normalized form. Group categorizations associated with the largest 25 (of 27) outliers in Figure 6 which are distinct from the main identified groups. PEEP is approximated in normalized pressure waveforms and indicated by dashed lines. Some outliers appear to be artifactual (from the data or estimation under stationarity). Others may be unique characterizations corresponding to extreme cases of VD, effects of patient posture, or heterogeneous breaths occurring under uncommonly used ventilator modes (e.g., spontaneous breathing present in 3.4% of breaths)

549 *Appendix A.3. Qualitative equivalence of labels via tSNE & UMAP*

550 Methodological choices may bias the segmentation process of LVS descriptors. The feature dimensional
551 reduction method used prior to DBSCAN labeling is strongly influential on the labeling process. Cluster
552 labels are qualitatively the same in nearly all cases for under application of tSNE and UMAP (Figs. A.9
553 and A.10). However, extracted characterizations for populous groupings may differ due to the geometries
554 of embedded points. Characterization of tSNE-oriented labels appear to be more representative of realized
555 breaths: tSNE projection of features tend to be more convex, which results in mean and median points lying
556 closer to realized data. [[What i'm trying to say here: UMAP coordinates can be more asymmetric and less
ball-like with tentacles, and loss of convexity means the 'center' can lie farther from the actual features.]]

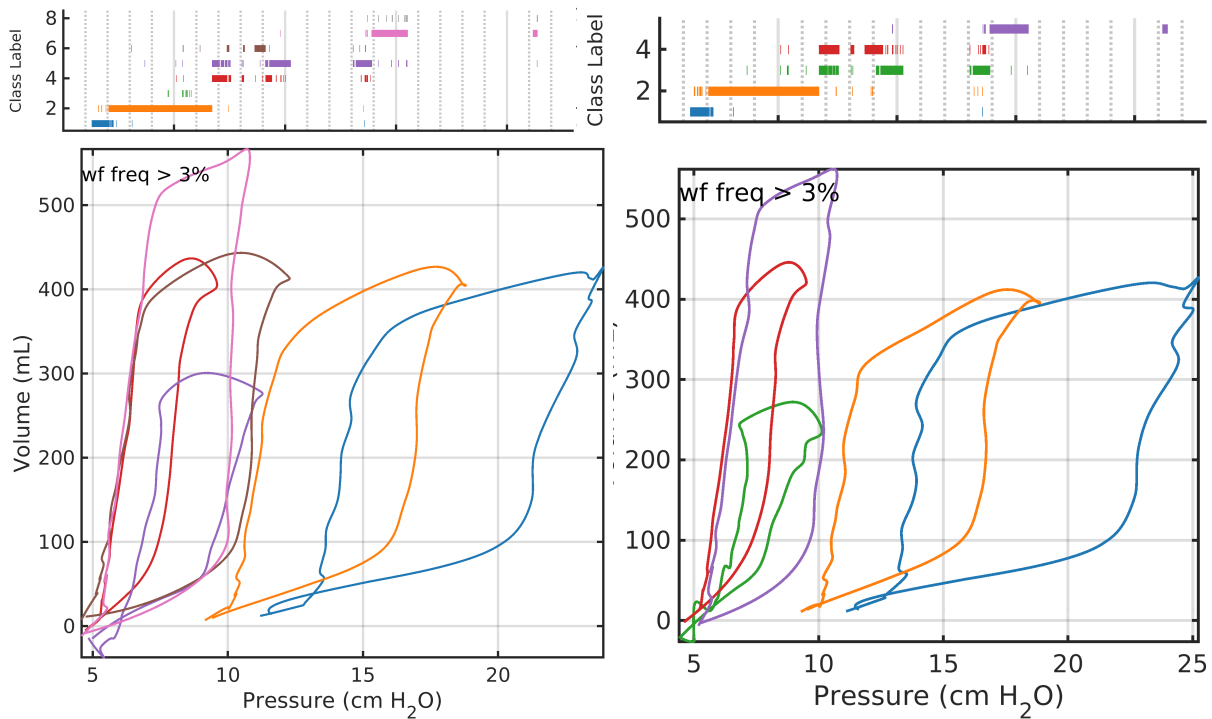


Figure A.9: Patient 101 clustering using tSNE (left) and UMAP (right) feature reduction stages. Identified phenotypes show qualitatively similar evolution although the tSNE-based characterization are more representative due poor representation of non-convex UMAP groupings by the component-wise median.

557

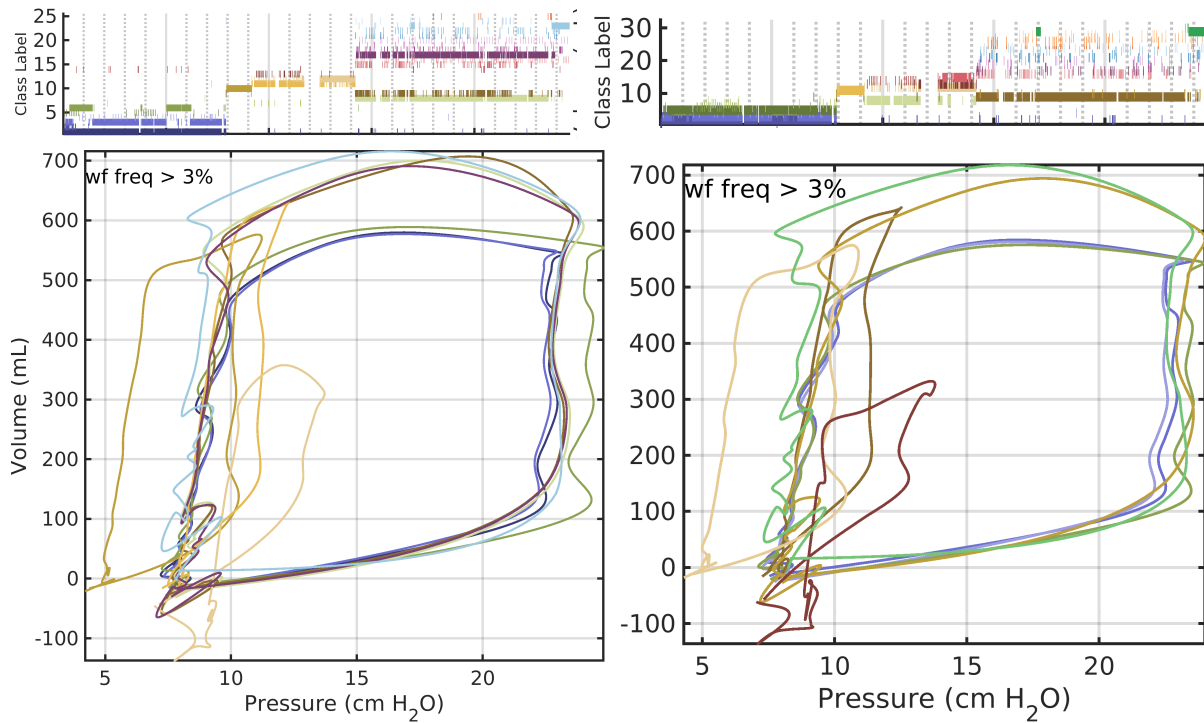


Figure A.10: Patient 124 clustering, as above

558 **Appendix B. Sample size vs. Sample description**

559 Broadly, waveform digitization transforms high-frequency temporal sampling of state processes into a
 560 lower-frequency, distributionally-descriptive form. This reduces the effective size of the problem while making
 561 it more dense. For a classification problem involving T samples of M -dimensional observations stored in
 562 an array $D \in \mathbb{R}^{T \times M}$, methods involving kernel or covariance processes require then calculating a matrix of
 563 dimension $M \times (T \times T) \times M$ in observation space or $T \times (M \times M) \times T$ in sample space. Decreasing the order
 564 of T and increasing that of M by a factor α benefits computational efficiency by replacing $D \in \mathbb{R}^{T \times M}$ with
 565 $\tilde{D} \in \mathbb{R}^{(T/\alpha) \times (\alpha M)}$. Specifically, calculating the observation covariance from \tilde{D} requires α^2 more storage but
 566 involves α^{-2} fewer calculations over the samples: $(\alpha M) \times (\alpha M)$ is calculated via $\alpha^{-1} T \times \alpha^{-1} T$
 567 rather than $T \times T$. Similarly, the summary sample space covariance of size $(T/\alpha) \times (T/\alpha)$ may be more
 568 dense than one built from un-summarized samples in $T \times T$, but it may be machine representable for larger
 569 values of T . Computational effects are important as $T \gg M$ in most practical applications, and additional
 570 statistical benefits arise from increasing the size of M .