

- 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*. 2<sup>nd</sup> 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*

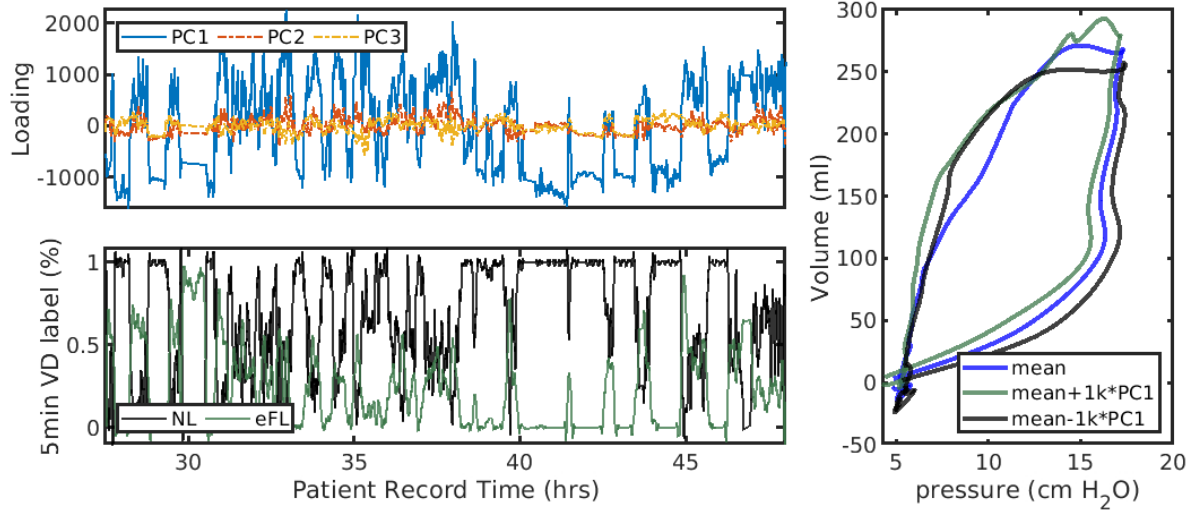


Figure A.7: The sign of PC1 loading roughly divides the VD classes in p111, label2. A threshold for the PC1 loading at zero roughly separates NL and eFL labels by 34%/65% and 85%/14%, respectively, with NL labels strongly associated with negative loadings. The optimal threshold ( $\sim 0.05$ ) offers 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  $\sim 3$  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*

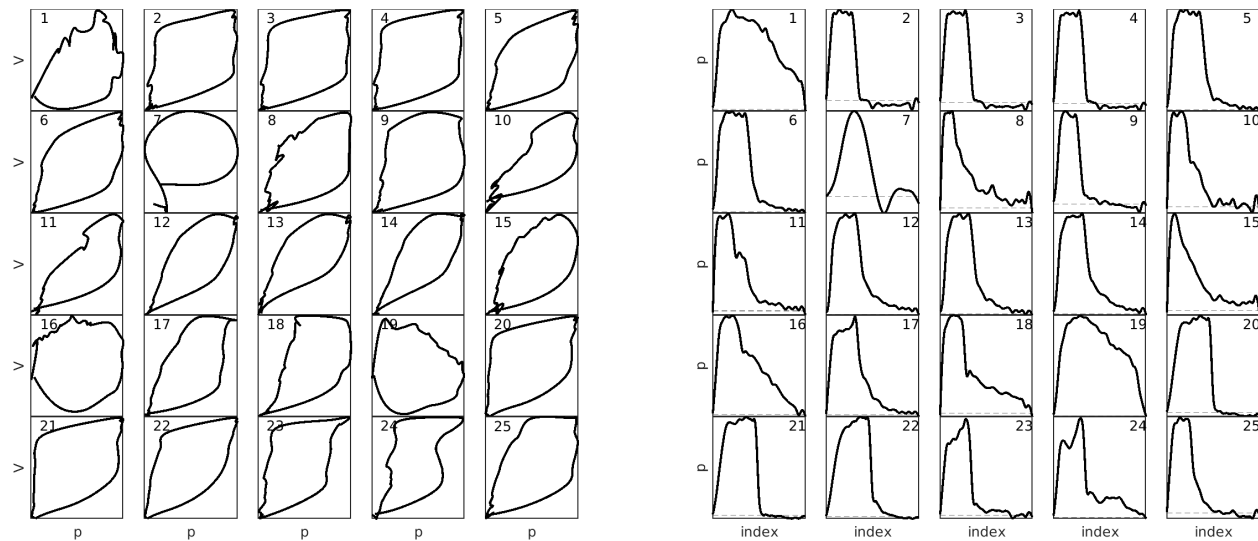


Figure A.8: Outlier pV (left) and pressure waveform (right) characterizations from two-stage cohort LVS phenotyping, shown in normalized form. Group categorizations associated with the largest 25 (of 27) outliers in Figure 6 which are distinct from the main identified groups. PEEP is approximated in 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.]]

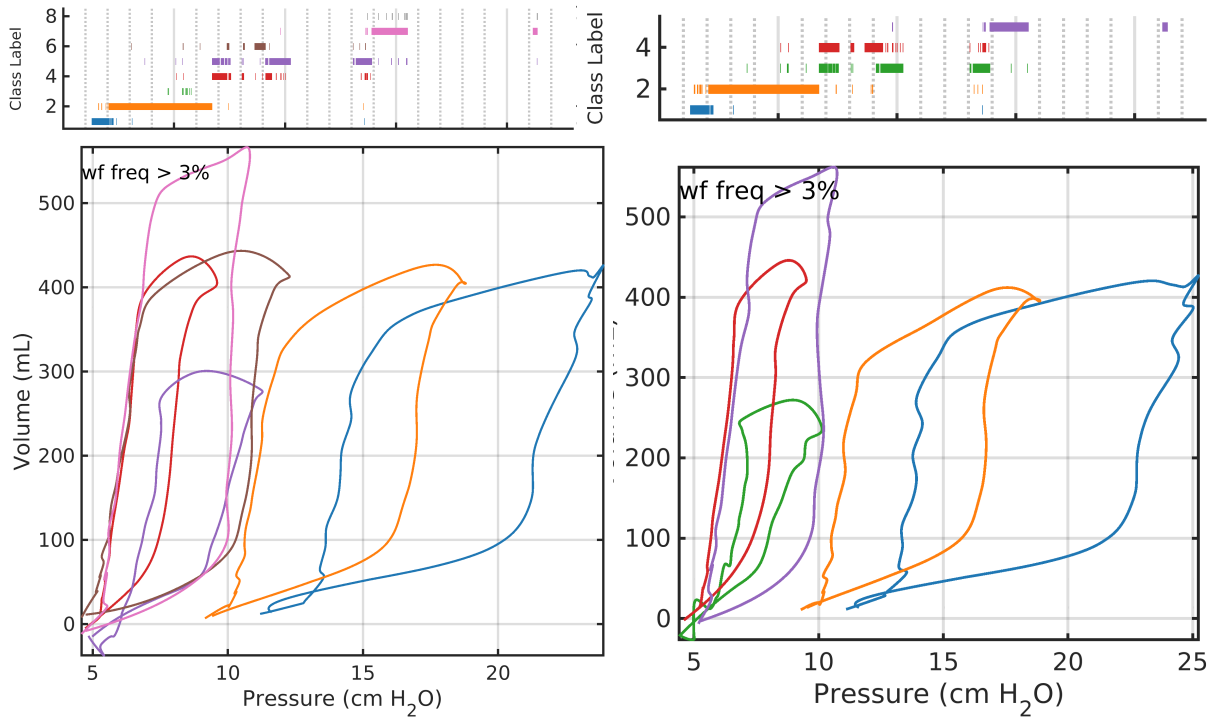


Figure A.9: Patient 101 clustering using tSNE (left) and UMAP (right) feature reduction stages. Identified phenotypes show qualitatively similar evolution although the tSNE-based characterization are more representative due poor representation of non-convex UMAP groupings by the component-wise median.

557

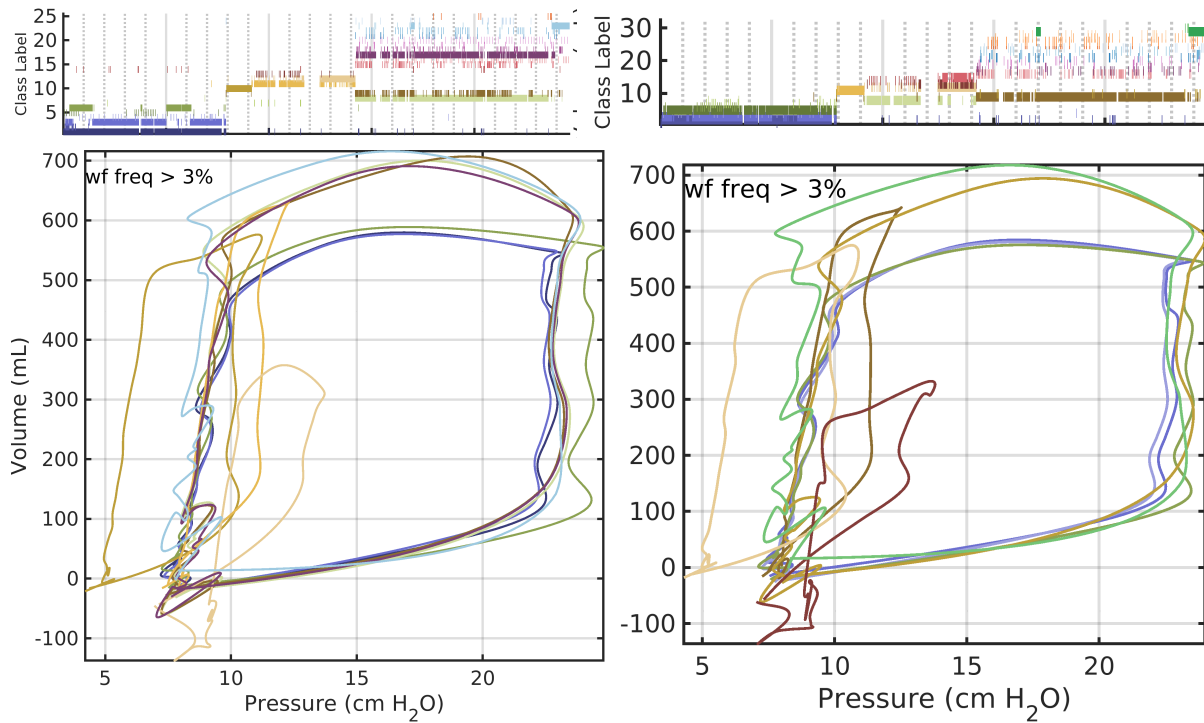


Figure A.10: Patient 124 clustering, as above

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  $\alpha$  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  $\alpha^2$  more storage but  
 566 involves  $\alpha^{-2}$  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  $M$ .