

# Geophysical Research Letters<sup>®</sup>

## RESEARCH LETTER

10.1029/2021GL096040

### Key Points:

- We develop a data-driven method that evaluates a velocity model using the K-means clustering and Rayleigh wave phase velocity dispersion
- The model evaluation method is applied to community velocity models, CVM-S4.26 and CVM-H15.1, in Southern California
- The result suggests that CVM-S4.26 gets an evaluation score ~3 times higher than that of CVM-H15.1 for structures in the top ~20 km

### Supporting Information:

Supporting Information may be found in the online version of this article.

### Correspondence to:

H. Qiu,  
[qiu@honrui@gmail.com](mailto:qiu@honrui@gmail.com);  
[hongruiq@mit.edu](mailto:hongruiq@mit.edu)

### Citation:

Xiong, N., Qiu, H., & Niu, F. (2021). Data-driven velocity model evaluation using K-means clustering. *Geophysical Research Letters*, 48, e2021GL096040. <https://doi.org/10.1029/2021GL096040>

Received 7 SEP 2021

Accepted 21 NOV 2021

### Author Contributions:

**Conceptualization:** Neng Xiong, Hongrui Qiu, Fenglin Niu  
**Data curation:** Hongrui Qiu  
**Formal analysis:** Neng Xiong  
**Methodology:** Neng Xiong, Hongrui Qiu  
**Supervision:** Hongrui Qiu, Fenglin Niu  
**Writing – original draft:** Neng Xiong  
**Writing – review & editing:** Neng Xiong, Hongrui Qiu, Fenglin Niu

## Data-Driven Velocity Model Evaluation Using K-Means Clustering

Neng Xiong<sup>1</sup> , Hongrui Qiu<sup>1,2</sup> , and Fenglin Niu<sup>1,3</sup> 

<sup>1</sup>Department of Earth, Environmental, and Planetary Sciences, Rice University, Houston, TX, USA, <sup>2</sup>Now at Department of Earth, Atmospheric, and Planetary Sciences, Massachusetts Institute of Technology, Cambridge, MA, USA, <sup>3</sup>State Key Laboratory of Petroleum Resources and Prospecting, and Unconventional Petroleum Research Institute, China University of Petroleum at Beijing, Beijing, China

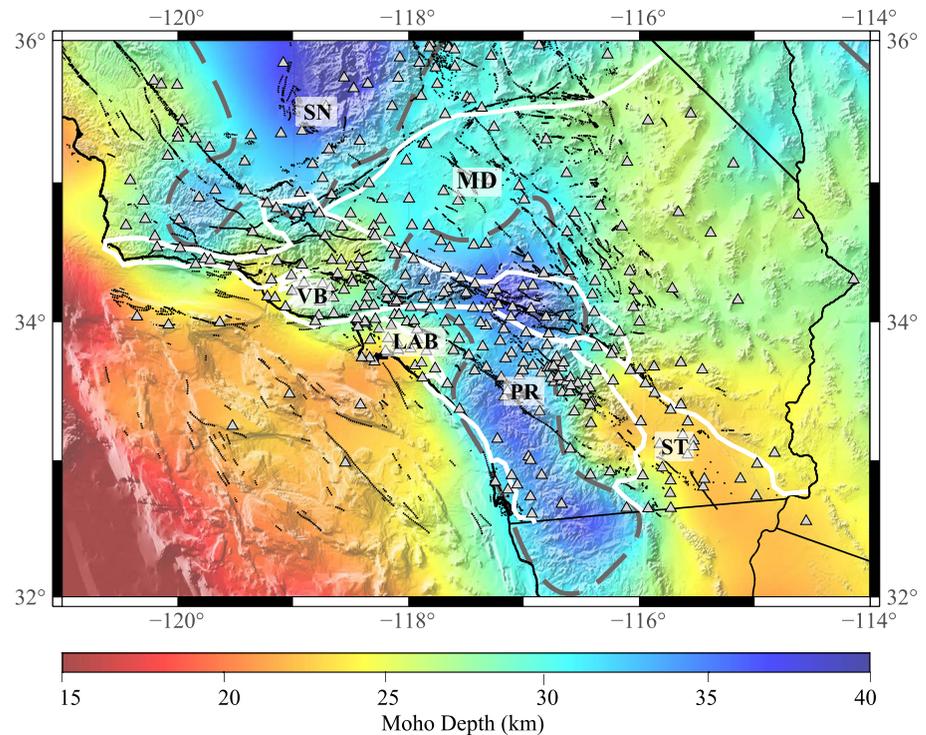
**Abstract** We develop a data-driven clustering method to evaluate a velocity model using surface wave velocity dispersion. This is done by first computing theoretical dispersion curves for 1-D velocity profiles of all the grid locations and then splitting the resulting dispersion curves into a certain number of groups via the K-means clustering. The observed dispersion curves are also clustered following the same procedure and the velocity model is assessed by comparing the spatial patterns obtained for the observed and synthetic data sets. The method is applied to evaluate two community velocity models in southern California, CVM-S4.26 and CVM-H15.1, using phase velocity maps derived for 3–16 s Rayleigh waves. We found a good correlation in the spatial distribution of clusters between the result of CVM-S4.26 and that of the observed data, suggesting that the CVM-S4.26 fits the observed dispersion maps better than the CVM-H15.1 in terms of features extracted from the clustering analysis.

**Plain Language Summary** With increasing volume of recorded seismic data, various velocity models are often derived for the same region using different data sets and seismic networks with different spatial coverage and resolution. Therefore, evaluating all the existing velocity models in the overlapping region can provide crucial information to future development of tomographic models, such as constructing a standard model by merging all the velocity models. As a machine learning technique, clustering analysis has proven its ability to extract hidden grouping features from large unlabeled data sets. In this study, we develop a simple workflow that utilizes a specific (K-means) clustering method to evaluate the velocity model. Instead of applying the clustering method directly to the velocity model, we first calculate theoretical predictions for a certain measurable parameter (phase velocity of Rayleigh wave) using the input model and assess the model by comparing the clustering results obtained for the synthetic and observed data sets. The proposed model evaluation method is applied to the well-maintained community velocity models, CVM-H15.1 and CVM-S4.26, in Southern California. The result suggests that CVM-S4.26 is much better than CVM-H15.1 for structures in the top ~20 km.

## 1. Introduction

With the increasing volume of data recorded by regional and global seismic networks, seismic tomography has become an important and powerful tool for understanding earth interior structure in the past decades. Southern California (SC; Figure 1) is one of the most active and imaged plate boundary regions. Velocity models that cover various depth (from near-surface to upper mantle) and spatial ranges with different resolutions were derived for this area (e.g., Berg et al., 2018; Lee et al., 2014; Lin et al., 2013; Roux et al., 2016). This is done by using different types of data sets, such as surface waves (e.g., Zigone et al., 2015) and teleseismic body waves (e.g., Schmandt & Humphreys, 2010), and inversion schemes, for example, by fitting travel-time (e.g., Fang et al., 2016) and full-waveform (e.g., Tape et al., 2010). Among these velocity models, the community velocity models (CVMs), CVM-H15.1 (Shaw et al., 2015) and CVM-S4.26 (Lee et al., 2014), are well maintained and often used as the starting model in travel-time based tomography studies (e.g., Qiu et al., 2019; Share et al., 2019).

Although CVM-H15.1 and CVM-S4.26 are both constructed through full-waveform inversion, differences between the two models are obvious (e.g., Figures 2a, 2b, 2d and 2e). This inconsistency is related to the choice of input data set (e.g., frequency range, station coverage, and usage of ambient noise data), inversion parameters (e.g., regularization and smoothing), and the starting model. Although some large-scale features (e.g., major

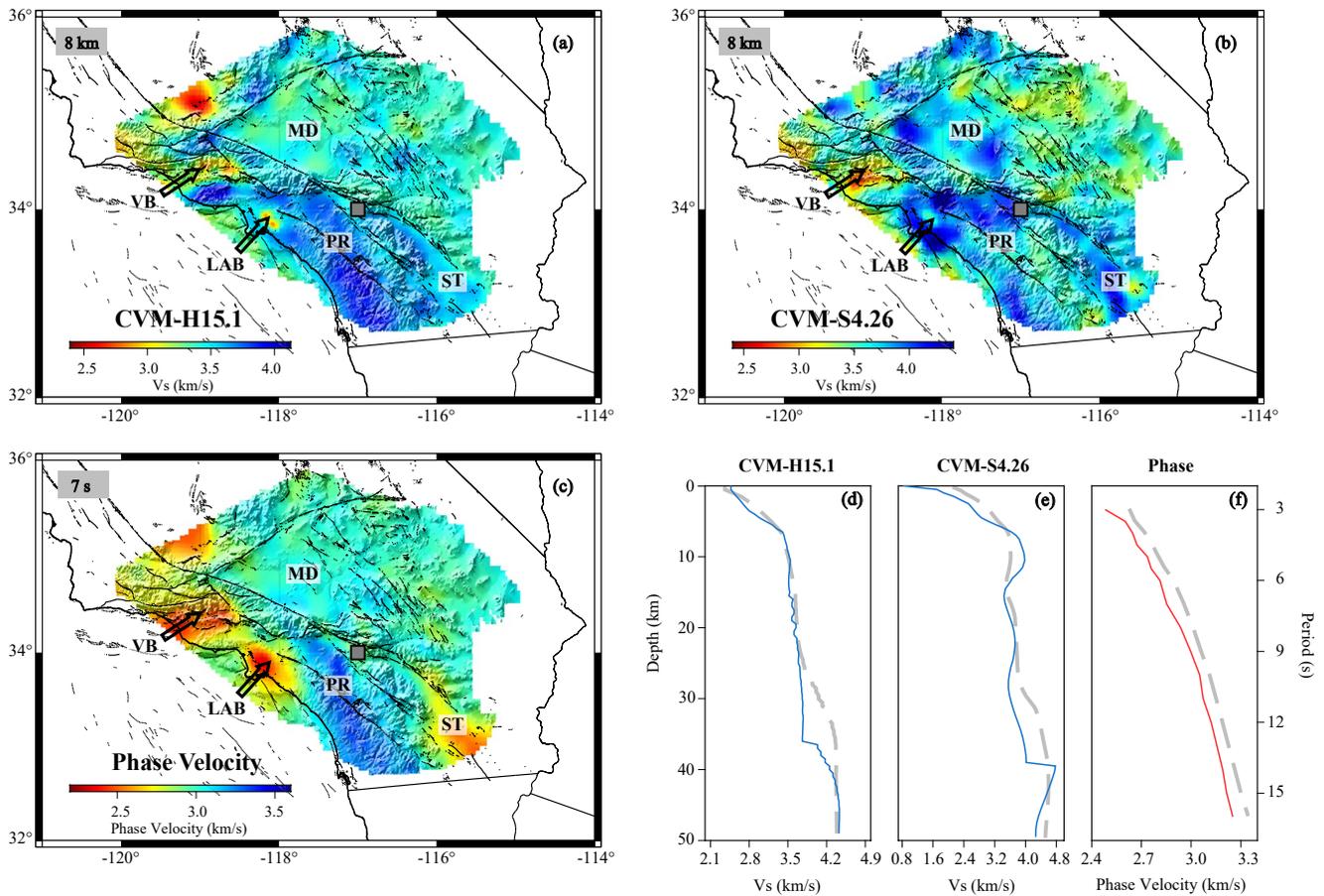


**Figure 1.** Map of Southern California plate boundary region. Background color indicates the Moho depth (Tape et al., 2012). Gray dashed line depicts the 31 km Moho depth contour. Gray triangles are the stations used in Qiu et al. (2019). White solid line outlines the boundaries of major geological provinces. MD: Mojave Desert, PR: Peninsular Ranges, LAB: LA Basin, VB: Ventura Basin, SN: Sierra Nevada, ST: Salton Trough.

geologic provinces; Figure 1) are seen in both models, it is still challenging to determine which model to use in studies that aim at improving or interpreting the velocity structure in SC. For instance, Qiu et al. (2019) demonstrated that the synthetic dispersion curves calculated using either CVM-H15.1 or CVM-S4.26 match poorly with the observed dispersion maps. However, through the 1-D  $V_s$  inversion (Herrmann, 2013) at each grid location, the misfit values are significantly reduced to a level comparable to the estimated uncertainties for both CVMs. This is likely due to the non-uniqueness of the inversion problem, which makes it difficult to evaluate which model is more realistic through the analysis of data misfit.

One way to assess the quality of a velocity model is through forward 3-D waveform simulation based on the wave equation. This is usually done by comparing earthquake recordings or empirical Green's functions retrieved from ambient noise data with synthetic waveforms simulated using the same source-receiver configuration (Imperator & Gallovič, 2017; Ma et al., 2008). However, the application of such a model validation method is limited by two main factors: (a) complicated evaluation scheme, that is, any inaccurate information in the velocity model along the ray path can contribute to the mismatch between synthetic and observed waveforms; and (b) intensive computational costs, particularly for observations at high frequencies (e.g., >1 Hz).

In recent years, machine learning has become more and more popular in extracting hidden features from large data sets in seismology (e.g., Bergen et al., 2019; Kong et al., 2018). Clustering analysis, as an unsupervised learning method, found success in mining different types of noise sources in continuous seismic recordings (Johnson et al., 2020; Snover et al., 2020). The nature of dividing data into groups with similar patterns makes clustering analysis suitable in dealing with large unlabeled data sets, such as seismic waveforms and velocity models. Eymold and Jordan (2019) applied the K-means clustering algorithm to the 1-D velocity profiles of CVM-S4.26 and discovered good correlation between surface geology features in SC and the resulting clustering pattern. However, it is important to note that, by directly clustering the 1-D velocity profiles, the obtained spatial pattern highly depends on the depth range of the input model (e.g., 0–50 km in Eymold & Jordan, 2019).



**Figure 2.** Vs maps at 8 km extracted from (a) CVM-H15.1 and (b) CVM-S4.26. (c) Phase velocity map at 7 s from Qiu et al. (2019). Gray square in (a)–(c) indicates the location of the vertical velocity profile shown in (d)–(f). Gray dashed line in (d)–(f) is the average profile of the entire study region.

Moreover, clustering results of the same region can also change with different input velocity models, and such difference is often hard to interpret, as the comparison does not involve data fitting to field measurements.

In this study, we propose a data-driven evaluation scheme for velocity models based on the K-means clustering method. This is done by first calculating synthetic surface wave velocity dispersion curves for all 1-D velocity profiles of an input velocity model, and then clustering the synthetic and observed velocity dispersion curves independently into a certain number of groups through clustering analysis. The velocity model is rated by estimating the similarity between spatial patterns obtained from the synthetic and observed dispersion data. The proposed method is applied to two velocity models in SC (CVM-H15.1 and CVM-S4.26). The two velocity models and the Rayleigh wave phase velocity dispersion maps measured by Qiu et al. (2019) that are used to assess the models are described in Section 2. The theoretical basis and workflow of the K-means clustering algorithm are reviewed and illustrated in Section 3. In Section 4, we show the spatial patterns of the clustering analysis for CVM-H15.1 and CVM-S4.26, and the evaluation of each model based on the observed phase velocity maps.

## 2. Data

The two CVMs, CVM-H15.1 and CVM-S4.26, analyzed in this study cover the SC plate boundary region (Figures 2a and 2b). Both models were extracted using the same grid size ( $0.05^\circ \times 0.05^\circ$ ) as the Rayleigh wave phase velocity dispersion maps. Depth of both CVMs were sampled with an interval of 500 m in the top 50 km. Except basin regions, the CVM-H15.1 was built upon an initial model derived from a local earthquake tomographic

inversion (Shaw et al., 2015). The model in basin areas is derived from more precise studies using borehole measurements and seismic reflection data (Süss & Shaw, 2003; Taborda et al., 2016), and held fixed during the wave-equation-based tomographic inversion. The CVM-H15.1 was derived utilizing data from 143 regional earthquakes recorded by 203 seismic stations (Shaw et al., 2015 and references therein). The CVM-S4.26 was constructed based on a different starting model via a similar tomographic inversion scheme (Lee et al., 2014) but using more earthquakes (160) and seismic stations (258). It is important to note that, in addition to earthquake data, ambient noise cross correlations calculated for pairs of stations are included in the inversion of CVM-S4.26.

The Rayleigh wave phase velocity dispersion maps used to evaluate the CVMs are discretized on a  $0.05^\circ \times 0.05^\circ$  grid and derived via Eikonal tomography from Qiu et al. (2019). The dispersion maps contain a total of 4,076 phase velocity dispersion curves ranging from 3 to 16 s. Figure shows the phase velocity map at the 7-s period. Clear phase velocity contrast can be seen across geologic provinces, such as high velocities in the Peninsular Ranges and low velocities in Salton Trough. To evaluate the CVMs via clustering analysis, the theoretical phase velocity dispersion curves are also calculated for all 1-D velocity profiles using the CPS package developed by Herrmann (2013). Both the observed and synthetic Rayleigh wave phase velocity dispersion curves are discretized into 17 data points from 3 to 16 s.

### 3. K-Means Clustering

In this study, we utilize the K-means clustering method to group a series of 1-D curves into a predetermined number of clusters. Let  $n$  be the number of the input 1-D curves, and  $K$  be the number of clusters. First,  $K$  1-D profiles are randomly chosen from the input data set as the initial centroids  $\{\mu_1, \mu_2, \mu_3, \dots, \mu_k\}$ . The Euclidean distances between each 1-D curve to all centroids are then calculated as the L2 norm between the two vectors:

$$D_k = \|x - \mu_k\|_2, \quad (1)$$

where  $x$  is the target velocity profile vector, and  $D$  is the distance vector that contains  $K$  number of values. Then, the target profile is assigned to its closest cluster, that is, the cluster yields the smallest distance. After all the profiles are assigned to a cluster, the centroid profile of each cluster is then updated as the average of all the profiles that belong to the cluster:

$$\mu'_k = \frac{\sum_{i=1}^{N_k} x_{ik}}{N_k}, \quad (2)$$

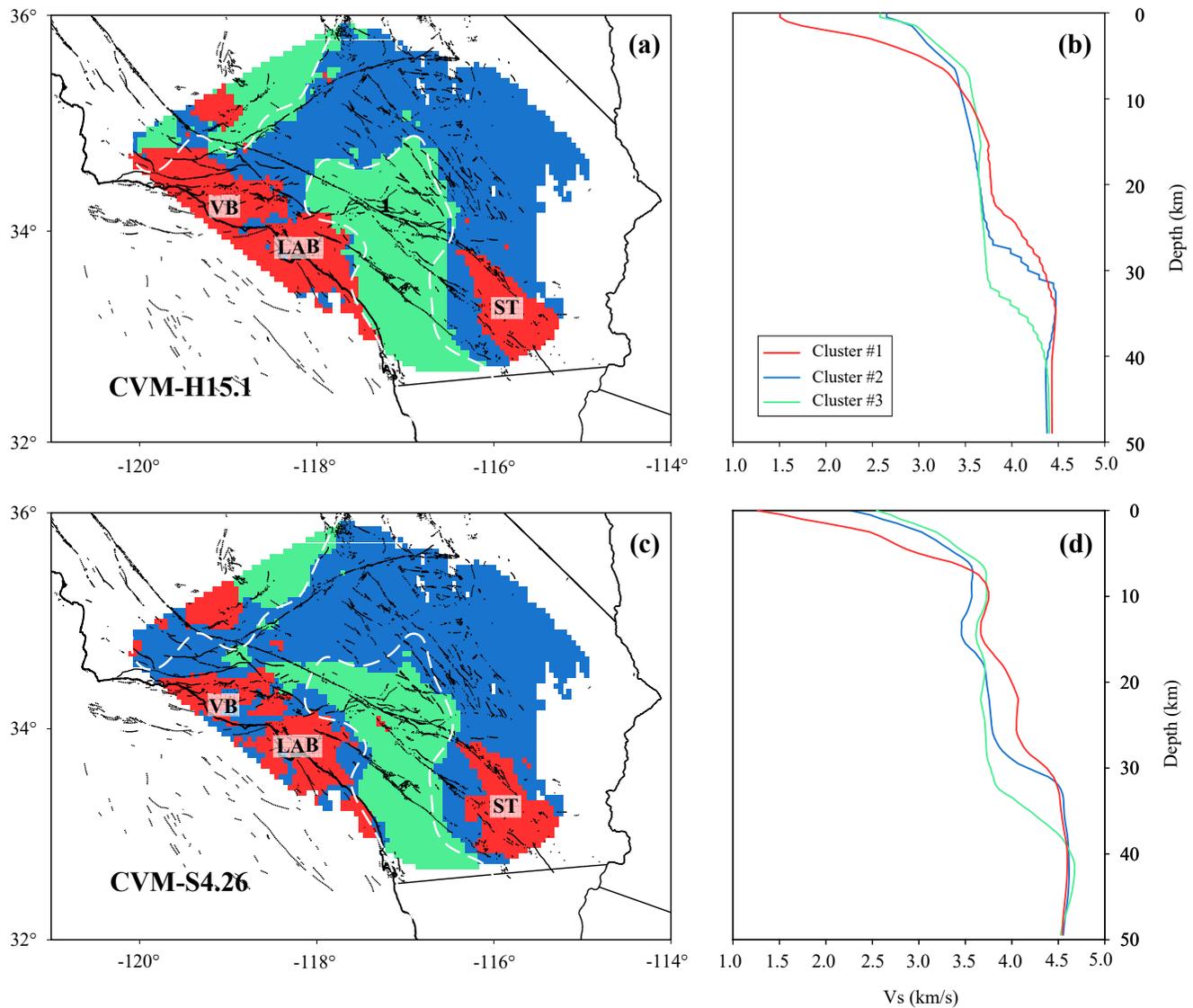
where  $N_k$  is the number of profiles in the  $k$ th cluster. If  $\mu'_k \neq \mu_k$  for any  $k$ th cluster, a new iteration of clustering process described by Equations 1 and 2 is performed, in which all data profiles are reassigned based on the updated centroids.

We note that the result of clustering analysis is sensitive to the choice of  $K$  value. The Elbow Method is often used to optimize the determination of  $K$  value (Eymold & Jordan, 2019). This is done by calculating the total distance of all data profiles to the corresponding centroid (of the cluster they assigned to), which is given by:

$$J(K) = \sum_{i=1}^n \sum_{k=1}^K \delta_i^k \|x_i - \mu_k\|^2 \quad (3a)$$

where

$$\delta_i^k = \begin{cases} 1, & \text{if } \min_j (|x_i - \mu_j|)^2 = |x_i - \mu_k|^2 \\ 0, & \text{otherwise} \end{cases} \quad (3b)$$



**Figure 3.**  $K = 3$  clustering results of (a) CVM-H15.1 and (c) CVM-S4.26. White dashed line is the 31 km Moho depth contour. Average Vs profile of each cluster of (b) CVM-H15.1 and (d) CVM-S4.26.

The optimal  $K$  value is determined as the knee of the objective function  $J(K)$ , where the gradient of the total variance flattens, indicating a diminishing return for increasing number of centroids. The clustering result may also be sensitive to the initial centroids if the objective function  $J(K)$  reaches a local minimum. Here, we run the clustering analysis 10 times using initial centroid locations generated randomly and keep the result with the lowest  $J(K)$ .

#### 4. Results

Figure 3 shows results of K-means clustering analysis performed directly on the Vs profiles of the CVM-H15.1 and CVM-S4.26 with an optimized  $K$  value of 3. This is similar to Eymold and Jordan (2019) but with the K-means clustering applied only to Vs in the top 50 km and grid cells covered by the phase velocity maps of Qiu et al. (2019). Similar large-scale spatial patterns can be seen for clustering results of both velocity models (Figures 3a and 3c). Cluster #1 (colored in red) covers regions with extremely low velocities at shallow depth, including sedimentary basins like Salton Trough and LA basin. For Clusters #2 (in blue) and #3 (in green), we overlay

the 31 km Moho depth contour resolved from Tape et al. (2012) onto the clustering maps (Figures 3a and 3c) and find a good correlation between the contour lines (white dashed curves) and boundaries between the two clusters.

The contour lines of the Moho interface at 31 km and the boundaries of Cluster #3 match particularly well for CVM-H15.1. In this case, the Moho depth variation dominates the clustering results (Figures 3b and 3d). This result is different from the more complicated pattern obtained in Eymold and Jordan (2019), which is likely because the 1-D Vp and Vs profiles at each grid cell are combined first before clustering and their study area is much larger. We note that the distribution of clusters could vary significantly if the depth range of the input Vs is changed, as the result would have no sensitivity to the Moho depth variation if structures only in the top 10 km are analyzed.

Although both models yield similar spatial patterns of the resulting clusters, obvious differences are still observed and difficult to interpret. In this study, however, we apply the clustering analysis to the synthetic phase dispersion curves calculated at all available grid cells for each CVM. Different from clustering of Vs profiles, the resulting spatial pattern of clusters from the synthetic phase velocity dispersion curves can be evaluated quantitatively using the observed phase velocity maps. Therefore, we first present the clustering analysis for the phase velocity maps derived by Qiu et al. (2019) and then evaluate each CVM by comparing the corresponding clustering result with that of the observed phase velocity maps.

#### 4.1. Clustering of the Observed Phase Velocity Maps

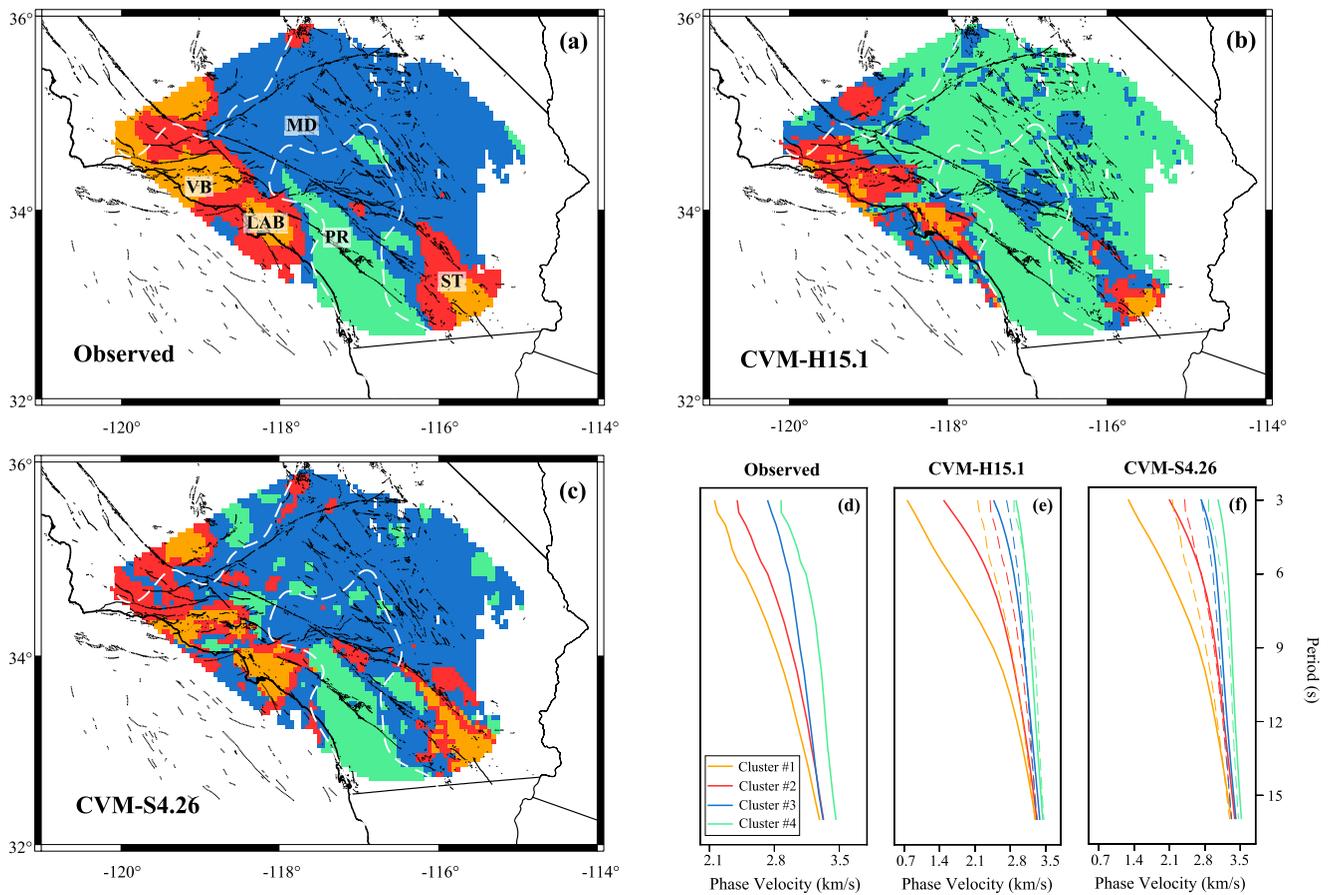
Figures 2c and 2f show a map-view and a 1-D profile of the Rayleigh wave phase velocity dispersion data obtained from Qiu et al. (2019), respectively. Compared to the Vs model (e.g., Figures 2d and 2e), the phase velocity profile (Figure 2f) is much smoother (e.g., no sharp velocity gradient due to Moho discontinuity) and sensitive to Vs values in a wide range of depth (Figure S1 in Supporting Information S1). Since the number of clusters  $K$  is a hyperparameter, we apply the Elbow Method and obtain the optimal  $K$  value as 4 (Figure S2 in Supporting Information S1). The clustering result of the observed phase velocity maps (Figure 4a) shows that Clusters #1 (in orange) and #2 (in red) mainly occupy the basin areas (e.g., LA basin and Salton Trough) with a relatively low phase velocity at short period (3–6 s), Cluster #3 (in blue) appears mostly in the Peninsular Ranges region, and Cluster #4 (in green) covers the Mojave Desert area.

#### 4.2. Clustering of Synthetic Phase Velocity Maps for CVM-H15.1

Similar to Section 4.1, we use  $K = 4$  in the clustering analysis of synthetic phase velocity dispersion curves calculated for CVM-H15.1 and the result is shown in Figure 4b. We note that, for a direct comparison, only dispersion curves calculated for grid cells covered by the data of Qiu et al. (2019) are included in the analysis. To ensure the colors assigned to clusters obtained for the CVM-H15.1 are consistent with those of the observed data, we use the centroid dispersion curve to label each cluster (Figures 4d and 4e). Although the resulting spatial pattern also highlights the Salton Trough, Los Angeles basin, and Ventura basin (i.e., low velocity anomalies at shallow depth; Figure 4e) with Clusters #1 and #2, the area is much smaller compared to those in Figure 4a. For each cluster label, we calculated the Jaccard index (Halkidi et al., 2002), the ratio between the sizes of intersection and union of two data sets, to estimate the similarity between two data sets and get the overall Jaccard index of 18.6% accounting for all clusters. We also compute the corresponding true positive rate (TPR) that is adopted in Eymold and Jordan (2019) for each cluster (Table S1 in Supporting Information S1).

#### 4.3. Clustering of Synthetic Phase Velocity Maps for CVM-S4.26

Clustering result of CVM-S4.26 using  $K = 4$  is shown in Figure 4c. A good spatial correlation is observed between Clusters #1 and #2 and basin areas. Moreover, the size of these two clusters agrees well with those in Figure 4a. Consistent with the clustering pattern for the observed phase velocity maps, the majority of the grid cells in Cluster #3 are also well confined within the Peninsular Ranges region (Figure 4c). Both Jaccard index and TPR for all clusters obtained from CVM-S4.26 are significantly higher than ( $\sim 2$ – $4$  times of) those of CVM-H15.1. More specifically, the overall Jaccard index of CVM-S4.26 is 57.4%, which is  $\sim 3$  times that of CVM-H15.1 (Table S1 in Supporting Information S1).



**Figure 4.**  $K = 4$  clustering result computed for (a) observed phase velocity and synthetic phase velocity of (b) CVM-H15.1 and (c) CVM-S4.26. Corresponding average phase velocity profile for each cluster (d)–(f). Dashed lines in (e) and (f) are average phase velocity profiles of each cluster shown in (d).

## 5. Discussion

In this study, we develop an alternative method to rate a velocity model via the K-means clustering method. This technique is applied to CVMs in SC using Rayleigh wave phase velocity maps derived from Qiu et al. (2019). Here, we further investigate the results by analyzing the K value, depth sensitivity kernel, and data misfit.

### 5.1. Selection of K Value

The K-means clustering analysis assigns similar data samples or profiles into the same cluster and is effective in extracting grouping features from large unlabeled data sets. However, the clustering result is dependent on the input number of clusters, that is, the K value. For clustering of 1-D Vs profiles (Figure 3) an optimal  $K = 3$  is applied, whereas an optimal  $K = 4$  is used in clustering of phase velocity dispersion curves in Section 4 (Figure 4). Since the effect of K value on the clustering result of 1-D Vs profiles is well discussed in Eymold and Jordan (2019), we focus on how the choice of K value alters the clustering of phase velocity dispersion curves and illustrate the results using  $K = 3$  in Figure S3 in Supporting Information S1 and  $K = 5$  in Figures S4–S6 in Supporting Information S1.

For  $K = 3$ , the number of extracted features from the clustering analysis is reduced compared to the case with  $K = 4$ . As expected, the spatial pattern shown in Figure S3a in Supporting Information S1 for the observed phase velocity dispersion curves is almost identical to that shown in Figure 4a after merging Clusters #1 (center of the basins) and #2 (edge of the basins) together. However, for the clustering of synthetic phase velocity dispersion curves (Figures S3b and S3c in Supporting Information S1), Cluster #3 in Figures 4a and 4c that primarily occupies the Peninsular Ranges region is missing from the  $K = 3$  results, indicating the difference between synthetic

dispersion curves in the center and at the edge of basins is much larger than the difference between basin and non-basin. This may be caused by the anomalously low phase velocities ( $<2$  km/s) in the period range of 3–5 s within LA basin, Ventura basin, and Salton Trough (red color areas in Figures S3b and S3c in Supporting Information S1), where significantly low  $V_s$  ( $<1$  km/s) at shallow depth are observed in both CVMs (Figures S3e and S3f in Supporting Information S1).

On the other hand, for  $K = 5$ , the clustering result of the observed phase velocity is highly dependent on the initialization, that is, the choice of starting centroids (Section 3). This is illustrated in Figure S4 in Supporting Information S1, where two different clustering patterns are obtained when two different starting centroids are randomly initialized. Such observed difference is greatly suppressed if we reduce the number of clusters from 5 to 4 by attributing the cluster in maroon to red and blue in Figures S4a and S4c in Supporting Information S1, respectively. The clustering result for the synthetic phase dispersion curves of CVMH-15.1 is also dependent on the centroid initialization (Figure S5 in Supporting Information S1), whereas the clustering result for CVM-S4.26 is less sensitive to the choice of starting centroids (Figure S6 in Supporting Information S1).

In conclusion, clustering results using  $K = 5$  are less stable than those of  $K = 3$  and  $K = 4$ , and the result of  $K = 3$  can be easily reproduced by merging two specific clusters obtained using  $K = 4$ . This likely suggests a maximum number of four dominating groups that can be extracted from the Rayleigh wave phase velocity dispersion curves between 3 and 16 s in the study area through clustering analysis, which justifies our choice of  $K = 4$  based on the Elbow Method result.

## 5.2. Depth Sensitivity

Figures 3a and 3c show the clustering results for 1-D  $V_s$  profiles in the top 50 km extracted from CVM-H15.1 and CVM-S4.26, respectively. The resulting spatial pattern yields two dominating structural features: basins (in red) with low velocities in the top 10 km and regions (in green) with a deep ( $>31$  km) Moho discontinuity. The clusters obtained using the observed phase velocity dispersion curves between 3 and 16 s, on the other hand, exhibit a different spatial pattern (Figures 4a and 4d). While the basins still stand out from the clustering results in Figure 4a, the other dominating structural feature outlined by the clustering analysis of dispersion curves is the Peninsular Ranges.

Considering Rayleigh wave phase velocities at periods  $<16$  s are most sensitive to structures in the top 20 km (Figure S1 in Supporting Information S1), the variation in Moho depth likely has little contribution to dispersion curves between 3 and 16 s. This is supported by the observation that the spatial pattern in Figure S7 in Supporting Information S1 derived using 1-D  $V_s$  profiles of CVM-S4.26 only in the top 20 km is consistent with that of the observed dispersion curves (Figure 4a). Therefore, we mainly evaluate the CVMs only in the top  $\sim 20$  km via clustering analysis of dispersion curves between 3 and 16 s. It is important to note that, in addition to extending the period range, we can also evaluate the velocity model at shallower depth by incorporating H/V ratio measurements from Berg et al. (2018).

## 5.3. Comparison With Data Misfit

We compute the data misfit of Rayleigh wave phase velocity for each CVM as the L2 norm between the observed and synthetic dispersion curves (Figure S8 in Supporting Information S1). The resulting misfit maps show similar patterns for both CVM-H15.1 and CVM-S4.26 with median values of  $\sim 0.5$  km/s. In general, basin regions yield large misfit values ( $>0.6$  km/s), while smaller values ( $<0.3$  km/s) are observed in the Mojave Desert and Peninsular Ranges. This suggests both models are similar in terms of fitting the phase velocity dispersion data. In contrast, our clustering-analysis-based evaluation method aims at comparing spatial patterns of the dominating structural features extracted independently from the observed and synthetic data sets, rather than focusing directly on the difference between them that is predominated by basin areas, and clearly shows that CVM-S4.26 is a better choice for structures in the top  $\sim 20$  km.

## 6. Conclusions

We develop, for the first time, a simple workflow to evaluate the velocity model via the K-means clustering method using observed surface wave phase velocity dispersion maps. This is done by first applying the K-means clustering analysis to synthetic phase velocity dispersion curves calculated for CVM-H15.1 and CVM-S4.26, and then validating each synthetic data set against the observed phase velocity maps obtained by Qiu et al. (2019). The resulting clustering pattern of both models is dominated by the distribution of sedimentary basins and major geologic provinces (e.g., Mojave Desert and Peninsular Ranges). Based on the comparison between clustering results of synthetic and observed dispersion curves, the Jaccard similarity coefficient averaged over all clusters is 57.4% for CVM-S4.26, which is more than three times as high as that of CVM-H15.1 (18.6%), suggesting the spatial pattern of clusters obtained from CVM-S4.26 matches much better with that of the observed data than CVM-H15.1. This is consistent with the fact that ambient noise cross correlation data is included in the inversion of CVM-S4.26 but not incorporated in the construction of CVM-H15.1.

Since the observed phase velocity maps between 3 and 16s are likely only sensitive to velocity structures in the top 20 km, other types of seismic data (e.g., H/V ratio, receiver function) that have higher sensitivity to a different depth range could be incorporated into the evaluation scheme to assess the part of velocity model at shallower or greater depth. The proposed clustering-based model evaluation method provides a simple and first-order rating system for any existing velocity models that complements the more sophisticated model validation studies based on 3-D full-waveform simulations and can provide crucial information to the future development of tomographic models, such as merging velocity models (e.g., determine the weighting of each velocity model in overlapping regions).

## Data Availability Statement

The Rayleigh wave phase velocity maps are obtained from Qiu et al. (2019) and accessible at <https://doi.org/10.17632/dt9x54dtrr.1>. The community velocity models were extracted using UCVMC (<https://github.com/SCECcode/UCVMC>). The Python module Scikit-Learn version 1.01 (Pedregosa et al., 2011) is used to perform the K-means clustering.

## References

- Berg, E. M., Lin, F.-C., Allam, A., Qiu, H., Shen, W., & Ben-Zion, Y. (2018). Tomography of Southern California via Bayesian joint inversion of Rayleigh wave ellipticity and phase velocity from ambient noise cross-correlations. *Journal of Geophysical Research: Solid Earth*, *123*, 9933–9949. <https://doi.org/10.1029/2018JB016269>
- Bergen, K. J., Johnson, P. A., de Hoop, M. V., & Beroza, G. C. (2019). Machine learning for data-driven discovery in solid Earth geoscience. *Science*, *363*(6433), eaau0323. <https://doi.org/10.1126/science.aau0323>
- Eymold, W. K., & Jordan, T. H. (2019). Tectonic regionalization of the Southern California crust from tomographic cluster analysis. *Journal of Geophysical Research: Solid Earth*, *124*(11), 11840–11865. <https://doi.org/10.1029/2019JB018423>
- Fang, H., Zhang, H., Yao, H., Allam, A., Zigone, D., Ben-Zion, Y., et al. (2016). A new algorithm for three-dimensional joint inversion of body wave and surface wave data and its application to the Southern California plate boundary region. *Journal of Geophysical Research: Solid Earth*, *121*, 3557–3569. <https://doi.org/10.1002/2015JB012702>
- Halkidi, M., Batistakis, Y., & Vazirgiannis, M. (2002). Cluster validity methods: Part I. *ACM SIGMOD Record*, *31*(2), 40–45. <https://doi.org/10.1145/565117.565124>
- Herrmann, R. B. (2013). Computer programs in seismology: An evolving tool for instruction and research. *Seismological Research Letters*, *84*(6), 1081–1088. <https://doi.org/10.1785/0220110096>
- Imperatori, W., & Gallovič, F. (2017). Validation of 3D velocity models using earthquakes with shallow slip: Case study of the 2014 Mw 6.0 South Napa, California, Event. *Bulletin of the Seismological Society of America*, *107*(2), 1019–1026. <https://doi.org/10.1785/0120160041>
- Johnson, C. W., Ben-Zion, Y., Meng, H., & Vernon, F. (2020). Identifying different classes of seismic noise signals using unsupervised learning. *Geophysical Research Letters*, *47*, e2020GL088353. <https://doi.org/10.1029/2020GL088353>
- Kong, Q., Trugman, D. T., Ross, Z. E., Bianco, M. J., Meade, B. J., & Gerstoft, P. (2018). Machine learning in seismology: Turning data into insights. *Seismological Research Letters*, *90*(1), 3–14. <https://doi.org/10.1785/0220180259>
- Lee, E. J., Chen, P., Jordan, T. H., Maechling, P. B., Denolle, M. A., & Beroza, G. C. (2014). Full-3-D tomography for crustal structure in southern California based on the scattering-integral and the adjoint-waveform methods. *Journal of Geophysical Research: Solid Earth*, *119*, 6421–6451. <https://doi.org/10.1002/2014JB011346>
- Lin, F.-C., Li, D., Clayton, R. W., & Hollis, D. (2013). High-resolution 3D shallow crustal structure in Long Beach, California: Application of ambient noise tomography on a dense seismic array. *Geophysics*, *78*, Q45–Q56. <https://doi.org/10.1190/geo2012-0453.1>
- Ma, S., Prieto, G. A., & Beroza, G. C. (2008). Testing community velocity models for Southern California using the ambient seismic field. *Bulletin of the Seismological Society of America*, *98*(6), 2694–2714. <https://doi.org/10.1785/0120080947>
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., et al. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, *12*(85), 2825–2830.
- Qiu, H., Lin, F.-C., & Ben-Zion, Y. (2019). Eikonal tomography of the Southern California plate boundary region. *Journal of Geophysical Research: Solid Earth*, *124*, 9755–9779. <https://doi.org/10.1029/2019JB017806>

## Acknowledgments

The authors thank W. K. Eymold for useful discussions and are grateful to all members of the Seismology and Tectonics at Rice University for comments and discussions. The authors thank the Editor, Dr. Daoyuan Sun, an anonymous reviewer and Dr. Weisen Shen for their constructive comments that help improve this paper. This study was supported by Rice University and the National Key Research and Development Program of China (No. 2017YFC1500300).

- Roux, P., Moreau, L., Lecointre, A., Hillers, G., Campillo, M., Ben-Zion, Y., et al. (2016). A methodological approach toward high-resolution seismic imaging of the San Jacinto Fault Zone using ambient noise recordings at a spatially-dense array. *Geophysical Journal International*, 206, 980–992. <https://doi.org/10.1093/gji/ggw193>
- Schmandt, B., & Humphreys, E. (2010). Seismic heterogeneity and small-scale convection in the southern California upper mantle. *Geochemistry, Geophysics, Geosystems*, 11, Q05004. <https://doi.org/10.1029/2010GC003042>
- Share, P. E., Guo, H., Thurber, C. H., Zhang, H., & Ben-Zion, Y. (2019). Seismic imaging of the Southern California Plate Boundary around the South-Central Transverse Ranges using double-difference tomography. *Pure and Applied Geophysics*, 176, 1117–1143. <https://doi.org/10.1007/s00024-018-2042-3>
- Shaw, J. H., Plesch, A., Tape, C., Süss, M. P., Jordan, T. H., Ely, G., et al. (2015). Unified structural representation of the southern California crust and upper mantle. *Earth and Planetary Science Letters*, 415, 1–15. <https://doi.org/10.1016/j.epsl.2015.01.016>
- Snover, D., Johnson, C. W., Bianco, M. J., & Gerstoft, P. (2020). Deep clustering to identify sources of urban seismic noise in Long Beach, California. *Seismological Research Letters*, 92, 1011–1022. <https://doi.org/10.1785/0220200164>
- Süss, M. P., & Shaw, J. H. (2003). P wave seismic velocity structure derived from sonic logs and industry reflection data in the Los Angeles basin, California. *Journal of Geophysical Research*, 108(B3), 2170. <https://doi.org/10.1029/2001JB001628>
- Taborda, R., Azizzadeh-Roodpish, S., Khoshnevis, N., & Cheng, K. (2016). Evaluation of the southern California seismic velocity models through simulation of recorded events. *Geophysical Journal International*, 205(3), 1342–1364. <https://doi.org/10.1093/gji/ggw085>
- Tape, C., Liu, Q., Maggi, A., & Tromp, J. (2010). Seismic tomography of the southern California crust based on spectral-element and adjoint methods. *Geophysical Journal International*, 180(1), 433–462. <https://doi.org/10.1111/j.1365-246X.2009.04429.x>
- Tape, C., Plesch, A., Shaw, J. H., & Gilbert, H. (2012). Estimating a continuous Moho surface for the California Unified Velocity Model. *Seismological Research Letters*, 83(4), 728–735. <https://doi.org/10.1785/0220110118>
- Zigone, D., Ben-Zion, Y., Campillo, M., & Roux, P. (2015). Seismic tomography of the Southern California plate boundary region from noise-based Rayleigh and Love waves. *Pure and Applied Geophysics*, 172(5), 1007–1032. <https://doi.org/10.1007/s00024-014-0872-1>