

Lawrence Berkeley National Laboratory

LBL Publications

Title

Data-Driven Velocity Model Evaluation Using K-Means Clustering

Permalink

<https://escholarship.org/uc/item/4bb730jp>

Journal

Geophysical Research Letters, 48(23)

ISSN

0094-8276

Authors

Xiong, Neng
Qiu, Hongrui
Niu, Fenglin

Publication Date

2021-12-16

DOI

10.1029/2021gl096040

Peer reviewed

1
2 **Data-driven Velocity Model Evaluation using K-means Clustering**
3

4 **Neng Xiong¹, Hongrui Qiu^{1,2*}, and Fenglin Niu¹**

5 ¹Department of Earth, Environmental and Planetary Sciences, Rice University, Houston, TX,
6 USA

7 ²Earth, Atmospheric and Planetary Sciences, Massachusetts Institute of Technology, Cambridge,
8 MA 02139, USA

9
10 Corresponding author: Hongrui Qiu (qiuhonrui@gmail.com; hongruiq@mit.edu)

11
12 **Key Points:**

- 13
- 14 • We develop a data-driven method that evaluates a velocity model using the K-means clustering and Rayleigh wave phase velocity dispersion
 - 15 • The model evaluation method is applied to community velocity models, CVM-S4.26 and
16 CVM-H15.1, in Southern California
 - 17 • The result suggests that CVM-S4.26 gets an evaluation score ~3 times higher than that of
18 CVM-H15.1 for structures in the top ~20 km
19

20 **Abstract**

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

32

33 **Plain Language Summary**

34 With increasing volume of recorded seismic data, various velocity models are often derived for
35 the same region using different datasets and seismic networks with different spatial coverage and
36 resolution. Therefore, evaluating all the existing velocity models in the overlapping region can
37 provide crucial information to future development of tomographic models, such as constructing a
38 standard model by merging all the velocity models. As a machine learning technique, clustering
39 analysis has proven its ability to extract hidden grouping features from large unlabeled datasets.
40 In this study, we develop a simple workflow that utilizes a specific (K-means) clustering method
41 to evaluate velocity model. Instead of applying the clustering method directly to the velocity
42 model, we first calculate theoretical predictions for a certain measurable parameter (phase
43 velocity of Rayleigh wave) using the input model and assess the model by comparing the
44 clustering results obtained for the synthetic and observed datasets. The proposed model
45 evaluation method is applied to the well-maintained community velocity models, CVM-H15.1
46 and CVMS-4.26, in Southern California. The result suggests that CVM-S4.26 is much better
47 than CVM-H15.1 for structures in the top ~20 km.

48 **1 Introduction**

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

61 Although CVM-H15.1 and CVM-S4.26 are both constructed through full-waveform
62 inversion, differences between the two models are obvious (e.g., Figs. 2a-b and 2d-e). This

63 inconsistency is related to the choice of input dataset (e.g., frequency range, station coverage,
64 usage of ambient noise data), inversion parameters (e.g., regularization, smoothing), and the
65 starting model. Although some large-scale features (e.g., major geologic provinces; Fig. 1) are
66 seen in both models, it is still challenging to determine which model to use in studies that aim at
67 improving or interpreting the velocity structure in SC. For instance, Qiu et al. (2019)
68 demonstrated that the synthetic dispersion curves calculated using either CVM-H15.1 or CVM-
69 S4.26 match poorly with the observed dispersion maps. However, through the 1-D Vs inversion
70 (Herrmann, 2013) at each grid location, the misfit values are significantly reduced to a level
71 comparable to the estimated uncertainties for both CVMs. This is likely due to the non-
72 uniqueness of the inversion problem, which makes it difficult to evaluate which model is more
73 realistic through the analysis of data misfit.

74 One way to assess the quality of a velocity model is through forward 3-D waveform
75 simulation based on the wave equation. This is usually done by comparing earthquake recordings
76 or empirical Green's functions retrieved from ambient noise data with synthetic waveforms
77 simulated using the same source-receiver configuration (Ma et al., 2008; Imperatori and
78 Gallovič, 2017). However, the application of such a model validation method is limited by two
79 main factors: (1) complicated evaluation scheme, i.e., any inaccurate information in the velocity
80 model along the ray path can contribute to the mismatch between synthetic and observed
81 waveforms; and (2) intensive computational costs, particularly for observations at high
82 frequencies (e.g., > 1 Hz).

83 In recent years, machine learning has become more and more popular in extracting hidden
84 features from large datasets in seismology (e.g., Bergen et al., 2019; Kong et al., 2018).
85 Clustering analysis, as an unsupervised learning method, found success in mining different types
86 of noise sources in continuous seismic recordings (Johnson et al., 2020; Snover et al., 2020). The
87 nature of dividing data into groups with similar pattern makes clustering analysis suitable in
88 dealing with large unlabeled datasets, such as seismic waveforms and velocity models. Eymold
89 and Jordan (2019) applied the K-means clustering algorithm to the 1-D velocity profiles of
90 CVM-S4.26 and discovered good correlation between surface geology features in SC and the
91 resulting clustering pattern. However, it is important to note that, by directly clustering the 1-D
92 velocity profiles, the obtained spatial pattern highly depends on the depth range of the input
93 model (e.g., 0-50 km in Eymold and Jordan, 2019). Moreover, clustering results of the same
94 region can also change with different input velocity models, and such difference is often hard to
95 interpret, as the comparison does not involve data fitting to field measurements.

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

108 2 Data

109 The two community velocity models, CVM-H15.1 and CVM-S4.26, analyzed in this study
 110 cover the SC plate boundary region (Figs. 2a-b). Both models were extracted using the same grid
 111 size ($0.05^\circ \times 0.05^\circ$) as the Rayleigh wave phase velocity dispersion maps. Depth of both CVMs
 112 were sampled with an interval of 500 m. Except basin regions, the CVM-H15.1 was built upon
 113 an initial model derived from a local earthquake tomographic inversion (Shaw et al., 2015). The
 114 model in basin areas is derived from more precise studies using borehole measurements and
 115 seismic reflection data (Süss and Shaw, 2003; Tabora et al, 2016), and held fixed during the
 116 wave-equation-based tomographic inversion. The CVM-H15.1 was derived utilizing data from
 117 143 regional earthquakes recorded by 203 seismic stations (Shaw et al., 2015 and references
 118 therein). The CVM-S4.26 was constructed based on a different starting model via a similar
 119 tomographic inversion scheme (Lee et al., 2014) but using more earthquakes (160) and seismic
 120 stations (258). It is important to note that, in addition to earthquake data, ambient noise cross
 121 correlations calculated for pairs of stations are included in the inversion of CVM-S4.26.

122 The Rayleigh wave phase velocity dispersion maps used to evaluate the CVMs are
 123 discretized on a $0.05^\circ \times 0.05^\circ$ grid and derived via Eikonal tomography from Qiu et al. (2019).
 124 The dispersion maps contain a total of 4076 phase velocity dispersion curves ranging from 3s to
 125 16s. Figure 2c shows the phase velocity map at 7s period. Clear phase velocity contrast can be
 126 seen across geologic provinces, such as high velocities in Peninsular Ranges and low velocities
 127 in Salton Trough. In order to evaluate the CVMs via clustering analysis, the theoretical phase
 128 velocity dispersion curves are also calculated for all 1-D velocity profiles using the CPS package
 129 developed by Herrmann (2013). Both the observed and synthetic Rayleigh wave phase velocity
 130 dispersion curves are discretized into 17 data points from 3s to 16s.

131 3 K-means clustering

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

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

137 where \mathbf{x} is the target velocity profile vector, and \mathbf{D} is the distance vector that contains K number
 138 of values. Then, the target profile is assigned to its closest cluster, i.e., the cluster yields the
 139 smallest distance. After all the profiles are assigned to a cluster, the centroid profile of each
 140 cluster is then updated as the average of all the profiles that belong to the cluster:

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

141 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
 142 iteration of clustering process described by equations (1) and (2) is performed, in which all data
 143 profiles are reassigned based on the updated centroids.

144 We note that the result of clustering analysis is sensitive to the choice of K value. The Elbow
 145 Method is often used to optimize the determination of K value (Eymold and Jordan 2019). This

146 is done by calculating the total distance of all data profiles to the corresponding centroid (of the
147 cluster they assigned to), which is given by:

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

148 where,

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

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

154 **4 Results**

155 Figure 3 shows results of K-means clustering analysis performed directly on the Vs profiles
156 of the CVM-H15.1 and CVM-S4.26 with an optimized K value of 3. This is similar to Eymold
157 and Jordan (2019) but with the K-means clustering applied only to Vs in the top 50 km and grid
158 cells covered by the phase velocity maps of Qiu et al. (2019). Similar large-scale spatial patterns
159 can be seen for clustering results of both velocity models (Figs. 3a and 3c). Cluster #1 (colored
160 in red) covers regions with extremely low velocities at shallow depth, including sedimentary
161 basins like Salton Trough and LA basin. For Clusters #2 (in blue) and #3 (in green), we overlay
162 the 31 km Moho depth contour resolved from Tape et al. (2012) onto the clustering maps (Figs.
163 3a and 3c) and find a good correlation between the contour lines (white dashed curves) and
164 boundaries between the two clusters.

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

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

181 **4.1. Clustering of the observed phase velocity maps**

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

192 **4.2. Clustering of synthetic phase velocity maps for CVM-H15.1**

193 Similar to section 4.1, we use $K = 4$ in the clustering analysis of synthetic phase velocity
 194 dispersion curves calculated for CVM-H15.1 and the result is shown in Figure 4b. We note that,
 195 for a direct comparison, only dispersion curves calculated for grid cells covered by the data of
 196 Qiu et al. (2019) are included in the analysis. To ensure the colors assigned to clusters obtained
 197 for the CVM-H15.1 are consistent with those of the observed data, we use the centroid
 198 dispersion curve to label each cluster (Figs. 4d and 4e). Although the resulting spatial pattern
 199 also highlights the Salton Trough, Los Angeles basin, and Ventura basin (i.e., low velocity
 200 anomalies at shallow depth; Fig. 4e) with Clusters #1 and #2, the area is much smaller compared
 201 to those in Figure 4a. For each cluster label, we calculated the Jaccard index (Halkidi et al.,
 202 2002), the ratio between the sizes of intersection and union of two datasets, to estimate the
 203 similarity between two datasets and get the overall Jaccard index of 18.6% accounting for all
 204 clusters. We also compute the corresponding true positive rate (TPR) that is adopted in Eymold
 205 & Jordan (2020) for each cluster (Table S1).

206 **4.3. Clustering of synthetic phase velocity maps for CVM-S4.26**

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

215 **5 Discussion**

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

220 **5.1. Selection of K value**

221 The K-means clustering analysis assigns similar data samples or profiles into the same
 222 cluster and is effective in extracting grouping features from large unlabeled datasets. However,
 223 the clustering result is dependent on the input number of clusters, i.e., the K value. In section 3,
 224 the optimal K is 3 for clustering of 1-D Vs profiles (Fig. 3), whereas an optimal $K = 4$ is used in

225 clustering of phase velocity dispersion curves in section 4 (Fig. 4). Since the effect of K value on
 226 the clustering result of 1-D Vs profiles is well discussed in Eymold and Jordan (2019), we focus
 227 on how the choice of K value alters the clustering of phase velocity dispersion curves and
 228 illustrate the results using $K = 3$ in Figure S3 and $K = 5$ in Figures S4-S6.

229 For $K = 3$, the number of extracted features from the clustering analysis is reduced compared
 230 to the case with $K = 4$. As expected, the spatial pattern shown in Fig. S3a for the observed phase
 231 velocity dispersion curves is almost identical to that shown in Fig. 4a after merging Clusters #1
 232 (center of the basins) and #2 (edge of the basins) together. However, for the clustering of
 233 synthetic phase velocity dispersion curves (Figs. S3b-c), Cluster #3 in Figs. 4a and 4c that
 234 primarily occupies the Peninsular Ranges region is missing from the $K = 3$ results, indicating the
 235 difference between synthetic dispersion curves in the center and at the edge of basins is much
 236 larger than the difference between basin and non-basin. This may be caused by the anomalously
 237 low phase velocities (< 2 km/s) in the period range of 3-5s within LA basin, Ventura basin, and
 238 Salton Trough (red color areas in Fig. S3b-c), where significantly low Vs (< 1 km/s) at shallow
 239 depth are observed in both CVMs (Figs. S3e-f).

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

248 In conclusion, clustering results using $K = 5$ are less stable than those of $K = 3$ and $K = 4$,
 249 and the result of $K = 3$ can be easily reproduced by merging two specific clusters obtained using
 250 $K = 4$. This likely suggests a maximum number of four dominating groups that can be extracted
 251 from the Rayleigh wave phase velocity dispersion curves between 3-16s in the study area
 252 through clustering analysis, which justifies our choice of $K = 4$ based on the Elbow Method
 253 result.

254 **5.2. Depth sensitivity**

255 Figures 3a and 3c show the clustering results for 1-D Vs profiles in the top 50 km extracted
 256 from CVM-H15.1 and CVM-S4.26, respectively. The resulting spatial pattern yields two
 257 dominating structural features: basins (in red) with low velocities in the top 10 km and regions
 258 (in green) with a deep (> 31 km) Moho discontinuity. The clusters obtained using the observed
 259 phase velocity dispersion curves between 3s and 16s, on the other hand, exhibit a different
 260 spatial pattern (Figs. 4a and 4d). While the basins still stand out from the clustering results in Fig.
 261 4a, the other dominating structural feature outlined by the clustering analysis of dispersion
 262 curves is the Peninsular Ranges.

263 Considering Rayleigh wave phase velocities at periods < 16 s are most sensitive to structures
 264 in the top 20 km (Fig. S1), the variation in Moho depth likely has little contribution to dispersion
 265 curves between 3s and 16s. This is supported by the observation that the spatial pattern in Fig. S7
 266 derived using 1-D Vs profiles of CVM-S4.26 only in the top 20 km is consistent with that of the
 267 observed dispersion curves (Fig. 4a). Therefore, we mainly evaluate the CVMs only in the top
 268 ~ 20 km via clustering analysis of dispersion curves between 3s and 16s. It is important to note

269 that, in addition to extending the period range, we can also evaluate the velocity model at
270 shallower depth by incorporating H/V ratio measurements from Berg et al. (2018).

271 **5.3. Comparison with data misfit**

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

282 **6 Conclusions**

283 We develop, for the first time, a simple workflow to evaluate velocity model via the K-means
284 clustering method using observed surface wave phase velocity dispersion maps. This is done by
285 first applying the K-means clustering analysis to synthetic phase velocity dispersion curves
286 calculated for CVM-H15.1 and CVM-S4.26, and then validating each synthetic dataset against
287 the observed phase velocity maps obtained by Qiu et al. (2019). The resulting clustering pattern
288 of both models is dominated by the distribution of sedimentary basins and major geologic
289 provinces (e.g., Mojave Desert and Peninsular Ranges). Based on the comparison between
290 clustering results of synthetic and observed dispersion curves, the Jaccard similarity coefficient
291 averaged over all clusters is 57.4% for CVM-S4.26, which is more than 3 times higher than that
292 of CVM-H15.1 (18.6%), suggesting the spatial pattern of clusters obtained from CVM-S4.26
293 matches much better with that of the observed data than CVM-H15.1. This is consistent with the
294 fact that ambient noise cross correlation data is included in the inversion of CVM-S4.26 but not
295 incorporated in the construction of CVM-H15.1.

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

305 **Data Availability Statement**

306 The Rayleigh wave phase velocity maps are obtained from Qiu et al. (2019) and accessible at
307 <https://doi.org/10.17632/dt9x54dtrr.1>. The community velocity models were extracted using
308 UCVMC (<https://github.com/SCECcode/UCVMC>). The Python module Scikit-Learn version
309 1.01 (Pedregosa et al., 2011) is used to perform the K-means clustering.

310

311 **Acknowledgements**

312 The authors thank W. K. Eymold for useful discussions and are grateful to all members of the
 313 Seismology and Tectonics at Rice University for comments and discussions. We thank the Editor
 314 Dr. Daoyuan Sun, an anonymous reviewer and Dr. Weisen Shen for their constructive comments
 315 that help improve this paper. This study was supported by Rice University.

316 **References**

- 317 Barak, S., Klemperer, S. L., & Lawrence, J. F. (2015). San Andreas Fault dip, Peninsular Ranges
 318 mafic lower crust and partial melt in the Salton Trough, Southern California, from ambient-
 319 noise tomography. *Geochemistry, Geophysics, Geosystems*, 16, 3946–3972.
 320 doi:10.1002/2015GC005970.
- 321 Berg, E. M., Lin, F.-C., Allam, A., Qiu, H., Shen, W., & Ben-Zion, Y. (2018). Tomography of
 322 Southern California via Bayesian joint inversion of Rayleigh wave ellipticity and phase
 323 velocity from ambient noise cross-correlations. *Journal of Geophysical Research: Solid
 324 Earth* 123, 9933–9949. <https://doi.org/10.1029/2018JB016269>
- 325 Bergen, K. J., Johnson, P. A., de Hoop, M. V., & Beroza, G. C. (2019). Machine learning for
 326 data-driven discovery in solid Earth geoscience. *Science*, 363(6433), eaau0323.
 327 <https://doi.org/10.1126/science.aau0323>
- 328 Eymold, W. K., & Jordan, T. H. (2019). Tectonic regionalization of the Southern California crust
 329 from tomographic cluster analysis. *Journal of Geophysical Research: Solid Earth*, 124,
 330 11,840–11,865. <https://doi.org/10.1029/2019JB018423>
- 331 Fang, H., Zhang, H., Yao, H., Allam, A., Zigone, D., Ben-Zion, Y., et al. (2016). A new
 332 algorithm for three-dimensional joint inversion of body wave and surface wave data and its
 333 application to the Southern California plate boundary region. *Journal of Geophysical
 334 Research: Solid Earth*, 121, 3557–3569. doi:10.1002/2015JB012702
- 335 Halkidi, M., Batistakis, Y., & Vazirgiannis, M. (2002). Cluster validity methods: part I. *ACM
 336 SIGMOD Record*, 31(2). <https://doi.org/10.1145/565117.565124>
- 337 Herrmann, R. B. (2013). Computer programs in seismology: An evolving tool for instruction and
 338 research. *Seismological Research Letters*, 84(6), 1081–1088.
 339 <https://doi.org/10.1785/0220110096>
- 340 Imperatori, W., & Gallovič, F. (2017). Validation of 3D Velocity Models Using Earthquakes
 341 with Shallow Slip: Case Study of the 2014 Mw 6.0 South Napa, California, Event. *Bulletin of
 342 the Seismological Society of America* 107 (2): 1019–1026. doi:
 343 <https://doi.org/10.1785/0120160041>
- 344 Johnson, C. W., Ben-Zion, Y., Meng, H., & Vernon, F. (2020). Identifying different classes of
 345 seismic noise signals using unsupervised learning. *Geophysical Research Letters*, 47,
 346 e2020GL088353. <https://doi.org/10.1029/2020GL088353>
- 347 Kong, Q., Trugman, D. T., Ross, Z. E., Bianco, M. J., Meade, B. J., & Gerstoft, P. (2018).
 348 Machine learning in seismology: Turning data into insights. *Seismological Research Letters*,
 349 90(1), 3–14. <https://doi.org/10.1785/0220180259>

- 350 Lee, E. J., Chen, P., Jordan, T. H., Maechling, P. B., Denolle, M. A., & Beroza, G. C. (2014).
351 Full-3-D tomography for crustal structure in southern California based on the scattering-
352 integral and the adjoint-waveform methods. *Journal of Geophysical Research: Solid Earth*,
353 119, 6421–6451. <https://doi.org/10.1002/2014JB011346>
- 354 Lin, F.-C., Moschetti, M. P., & Ritzwoller, M. H. (2008). Surface wave tomography of the
355 western United States from ambient seismic noise: Rayleigh and Love wave phase velocity
356 maps. *Geophysical Journal International*, 173(1), 281–298. <https://doi.org/10.1111/j.1365-246X.2008.03720.x>
- 357
- 358 Lin, F.-C., Li, D., Clayton, R. W., & Hollis, D. (2013). High-resolution 3D shallow crustal
359 structure in Long Beach, California: Application of ambient noise tomography on a dense
360 seismic array. *Geophysics*, 78: Q45-Q56. <https://doi.org/10.1190/geo2012-0453.1>
- 361 Ma, S., Prieto, G. A., & Beroza, G. C. (2008). Testing Community Velocity Models for Southern
362 California Using the Ambient Seismic Field. *Bulletin of the Seismological Society of America*,
363 98 (6): 2694–2714. doi: <https://doi.org/10.1785/0120080947>
- 364 Magistrale, H., McLaughlin, K., & Day, S. (1996). A geology-based 3D velocity model of the
365 Los Angeles basin sediments. *Bulletin of the Seismological Society of America*, 86(4): 1161–
366 1166.
- 367 Magistrale, H., Day, S., Clayton, R.W., & Graves, R. (2000). The SCEC southern California
368 reference three-dimensional seismic velocity model version 2. *Bulletin of the Seismological*
369 *Society of America*, 90(6B), S65–S76. <https://doi.org/10.1785/0120000510>
- 370 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., et al. (2011).
371 Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12(85),
372 2825–2830. Retrieved from <http://jmlr.org/papers/v12/pedregosa11a.html>
- 373 Qiu, H., Lin, F.-C., & Ben-Zion, Y. (2019). Eikonal Tomography of the Southern California
374 Plate Boundary Region. *Journal of Geophysical Research: Solid Earth*, 124, 9755–9779.
375 <https://doi.org/10.1029/2019JB017806>
- 376 Roux, P., Moreau, L., Lecointre, A., Hillers, G., Campillo, M., Ben-Zion, Y., Zigone, D., &
377 Vernon, F. (2016). A methodological approach toward high-resolution seismic imaging of
378 the San Jacinto Fault Zone using ambient noise recordings at a spatially-dense array.
379 *Geophysical Journal International*, 206, 980–992.
- 380 Schmandt, B., & Humphreys, E. (2010). Seismic heterogeneity and small-scale convection in the
381 southern California upper mantle. *Geochemistry, Geophysics, Geosystems*, 11, Q05004.
382 <https://doi.org/10.1029/2010GC003042>
- 383 Share, PE., Guo, H., Thurber, C. H., Zhang, H., & Ben-Zion, Y. (2019). Seismic Imaging of the
384 Southern California Plate Boundary around the South-Central Transverse Ranges Using
385 Double-Difference Tomography. *Pure and Applied Geophysics*, 176, 1117–1143.
386 <https://doi.org/10.1007/s00024-018-2042-3>
- 387 Shaw, J. H., Plesch, A., Tape, C., Süß, M. P., Jordan, T. H., Ely, G., et al. (2015). Unified
388 structural representation of the southern California crust and upper mantle. *Earth and*
389 *Planetary Science Letters*, 415, 1. <https://doi.org/10.1016/j.epsl.2015.01.016>

- 390 Snover, D., Johnson, C. W., Bianco, M. J., & Gerstoft, P. (2020). Deep Clustering to Identify
391 Sources of Urban Seismic Noise in Long Beach, California. *Seismological Research Letters*,
392 92, 1011–1022. <https://doi.org/10.1785/0220200164>.
- 393 Süss, M. P., & Shaw, J. H. (2003). P wave seismic velocity structure derived from sonic logs and
394 industry reflection data in the Los Angeles basin, California. *Journal of Geophysical*
395 *Research*, 108(B3), 2170. <https://doi.org/10.1029/2001JB001628>
- 396 Taborda, R., Azizzadeh-Roodpish, S., Khoshnevis, N., & Cheng, K. (2016). Evaluation of the
397 southern California seismic velocity models through simulation of recorded events.
398 *Geophysical Journal International*, 205(3), 1342–1364. <https://doi.org/10.1093/gji/ggw085>
- 399 Tape, C., Liu, Q., Maggi, A., & Tromp, J. (2010). Seismic tomography of the southern California
400 crust based on spectral-element and adjoint methods. *Geophysical Journal International*,
401 180(1): 433–462. <https://doi.org/10.1111/j.1365-246X.2009.04429.x>
- 402 Tape, C., Plesch, A., Shaw, J. H., & Gilbert, H. (2012). Estimating a Continuous Moho Surface
403 for the California Unified Velocity Model. *Seismological Research Letters*, 83(4): 728–735.
404 doi: <https://doi.org/10.1785/0220110118>
- 405 Yang, Y., & Forsyth, D. W. (2006). Rayleigh wave phase velocities, small-scale convection, and
406 azimuthal anisotropy beneath southern California. *Journal of Geophysical Research*, 111,
407 B07306. <https://doi.org/10.1029/2005JB004180>
- 408 Zigone, D., Ben-Zion, Y., Campillo, M., & Roux, P. (2015). Seismic tomography of the
409 Southern California plate boundary region from noise-based Rayleigh and Love waves. *Pure*
410 *and Applied Geophysics*, 172(5), 1007–1032. <https://doi.org/10.1007/s00024-014-0872-1>
411
412

413 **Figure caption**

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

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

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

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

Figure 1.

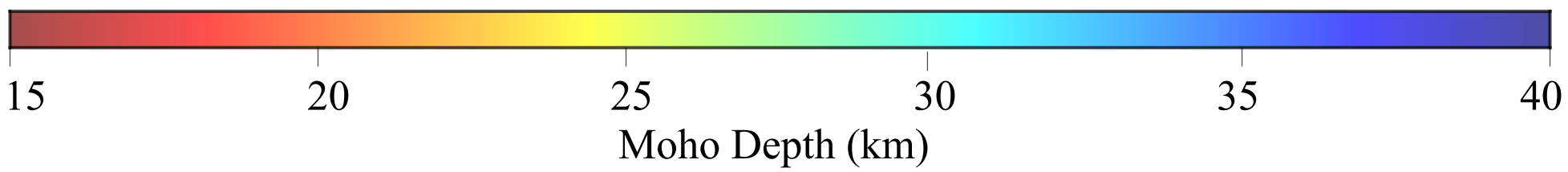
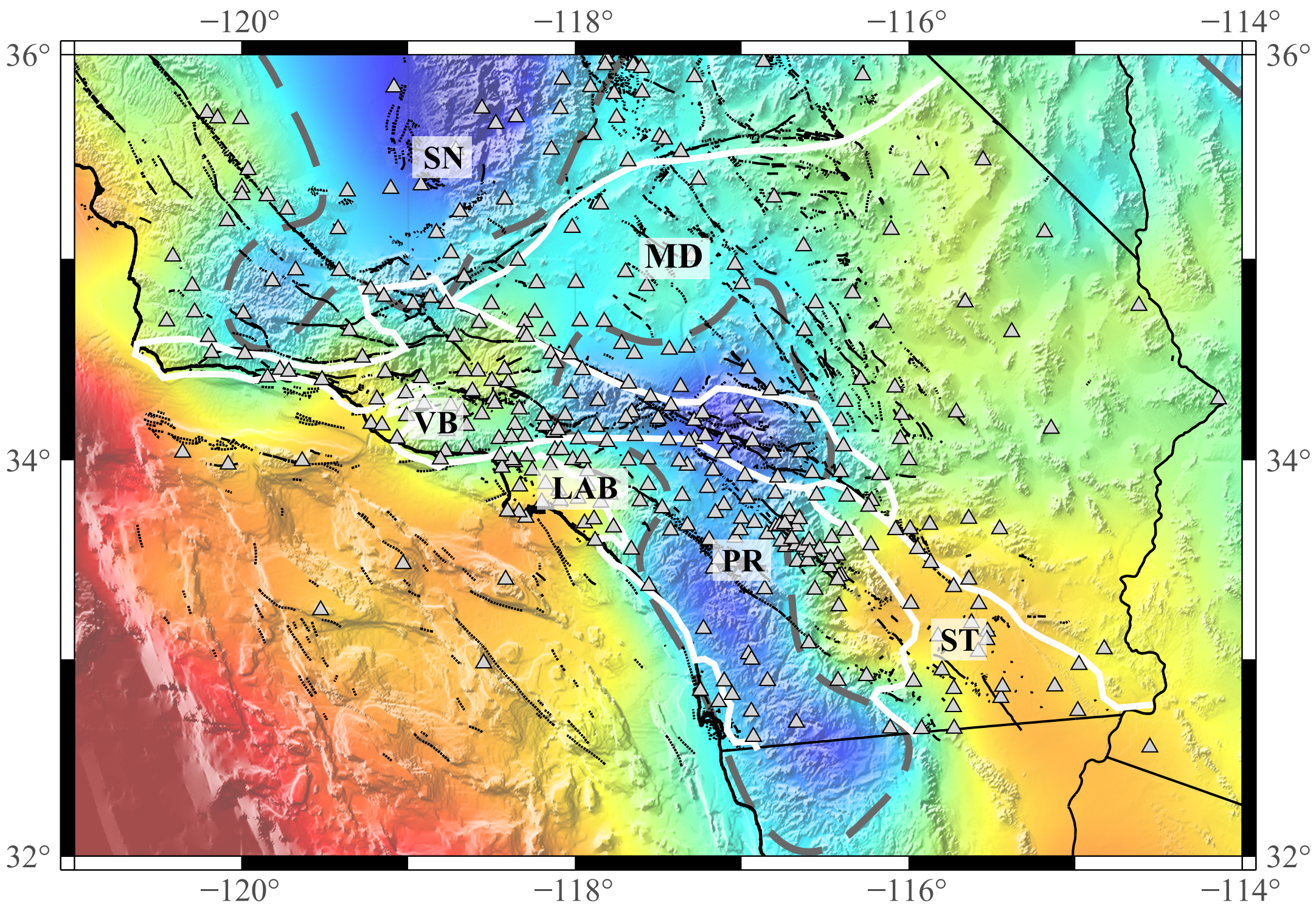


Figure 2.

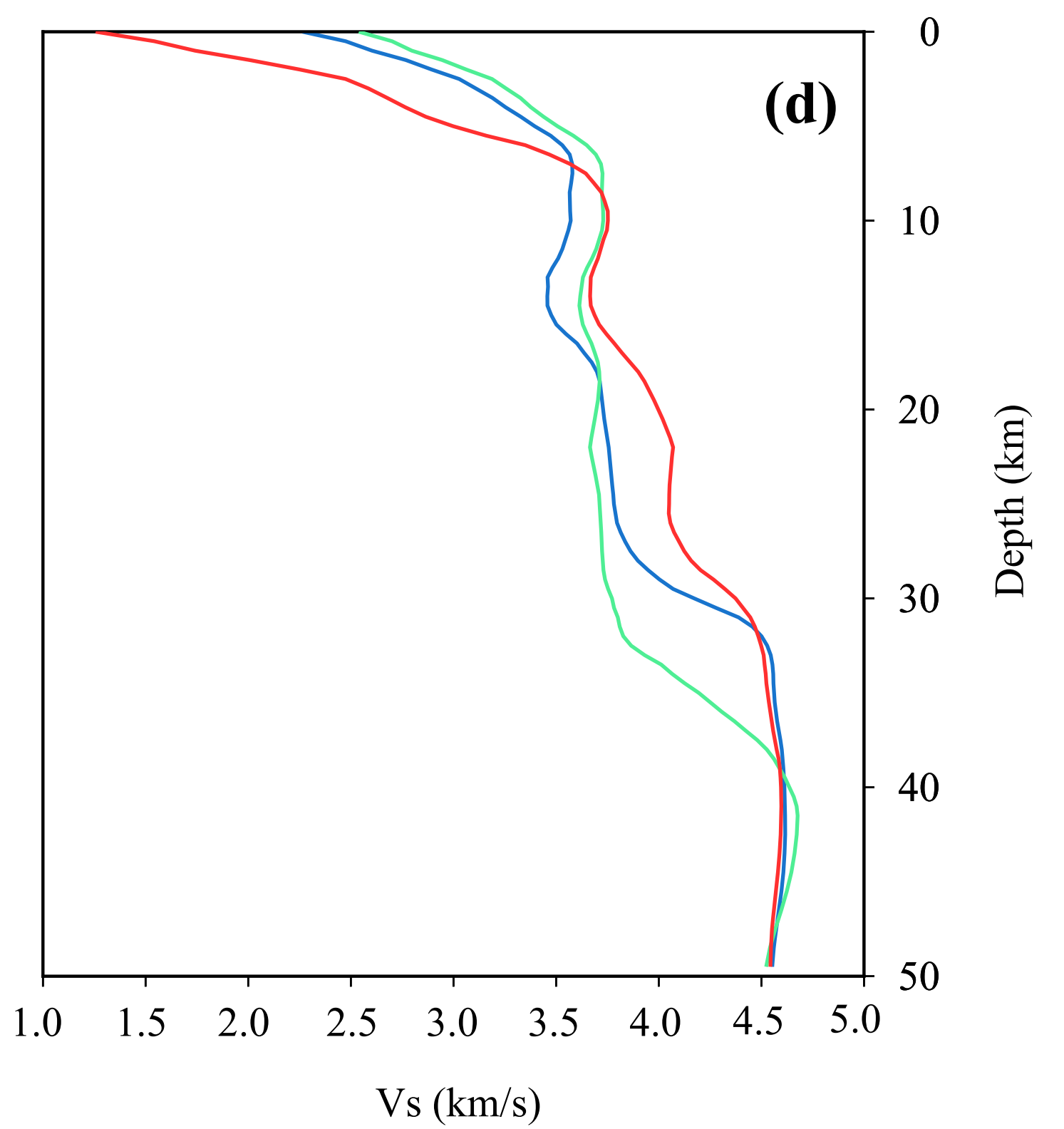
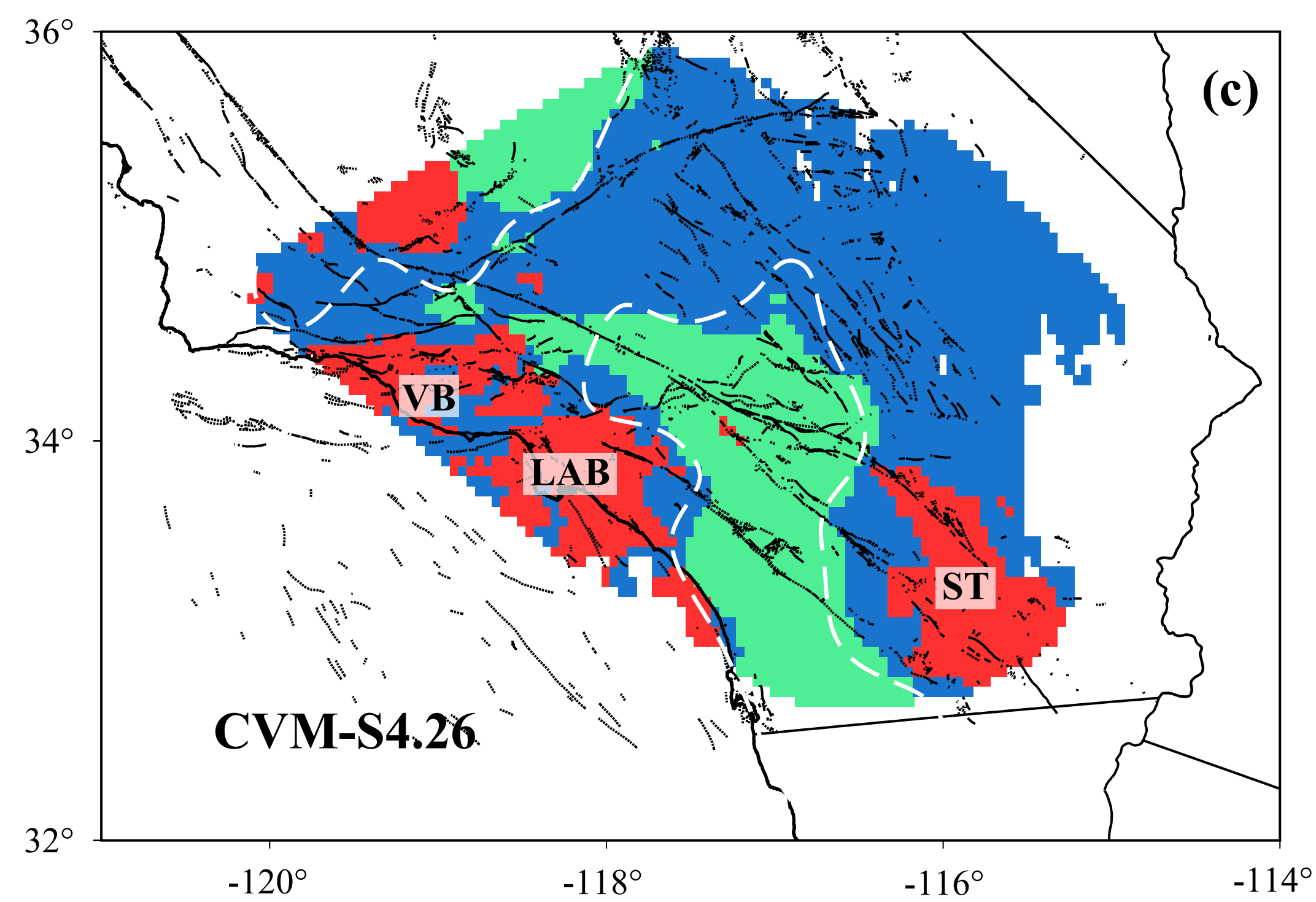
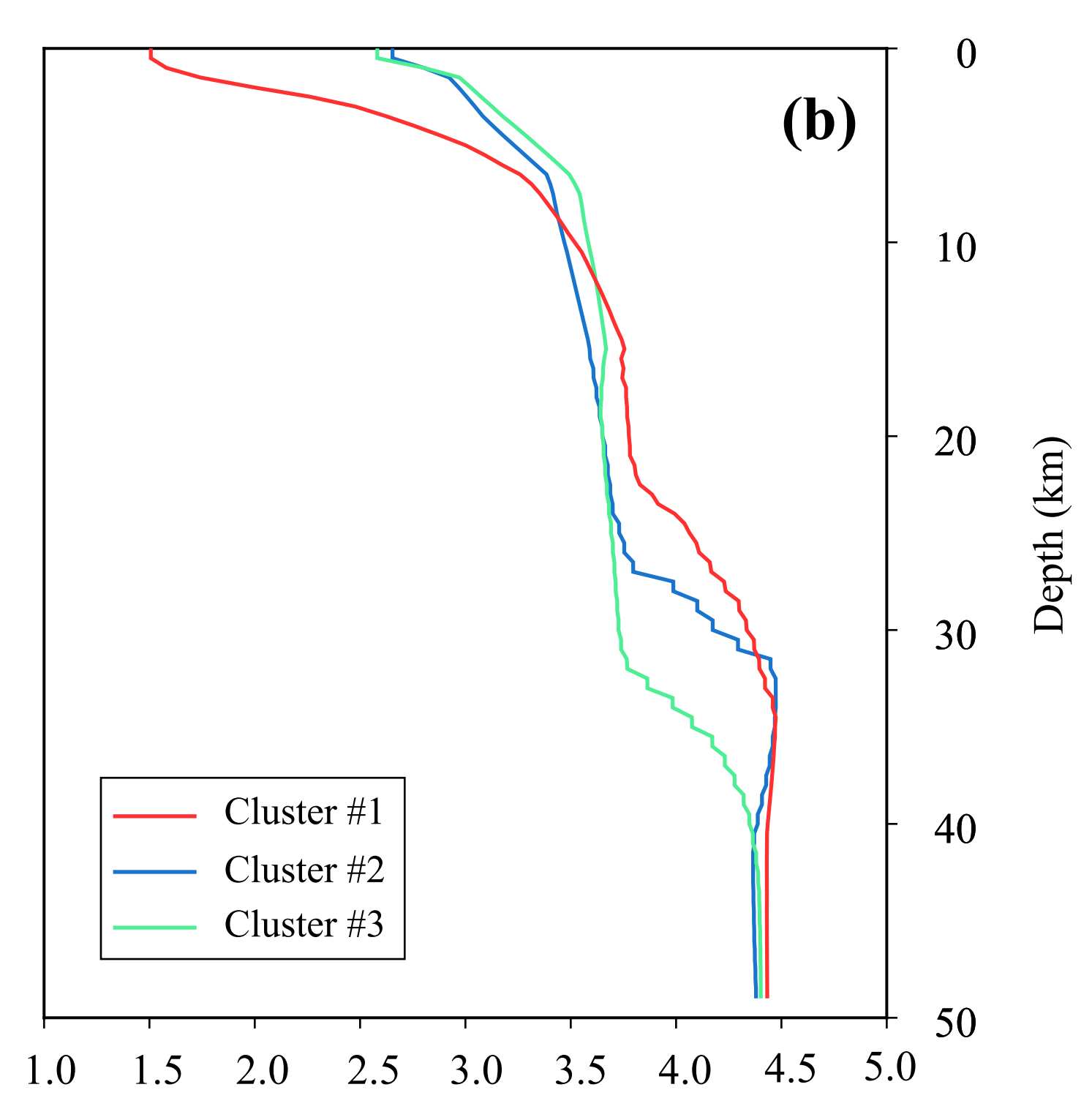
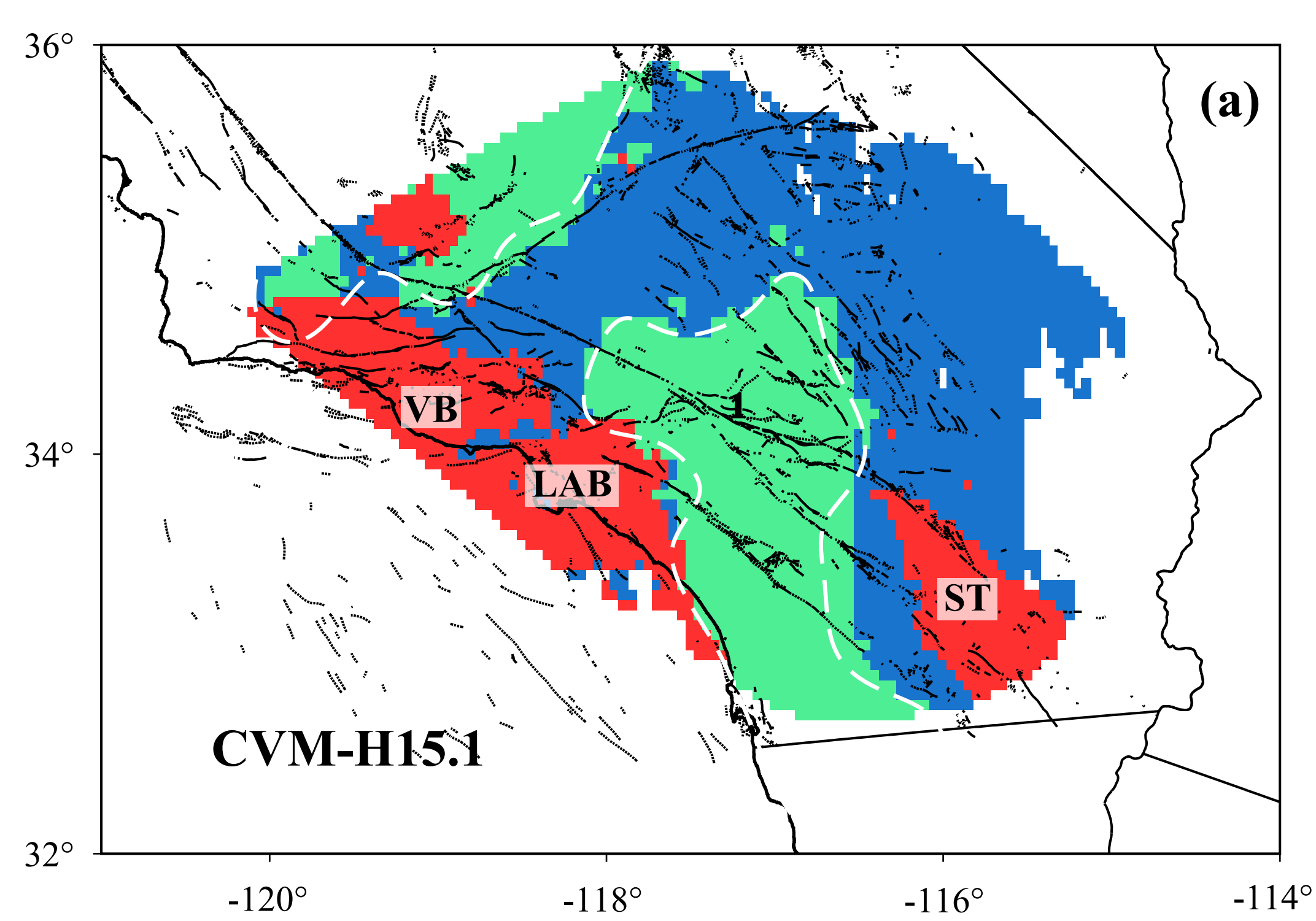


Figure 3.

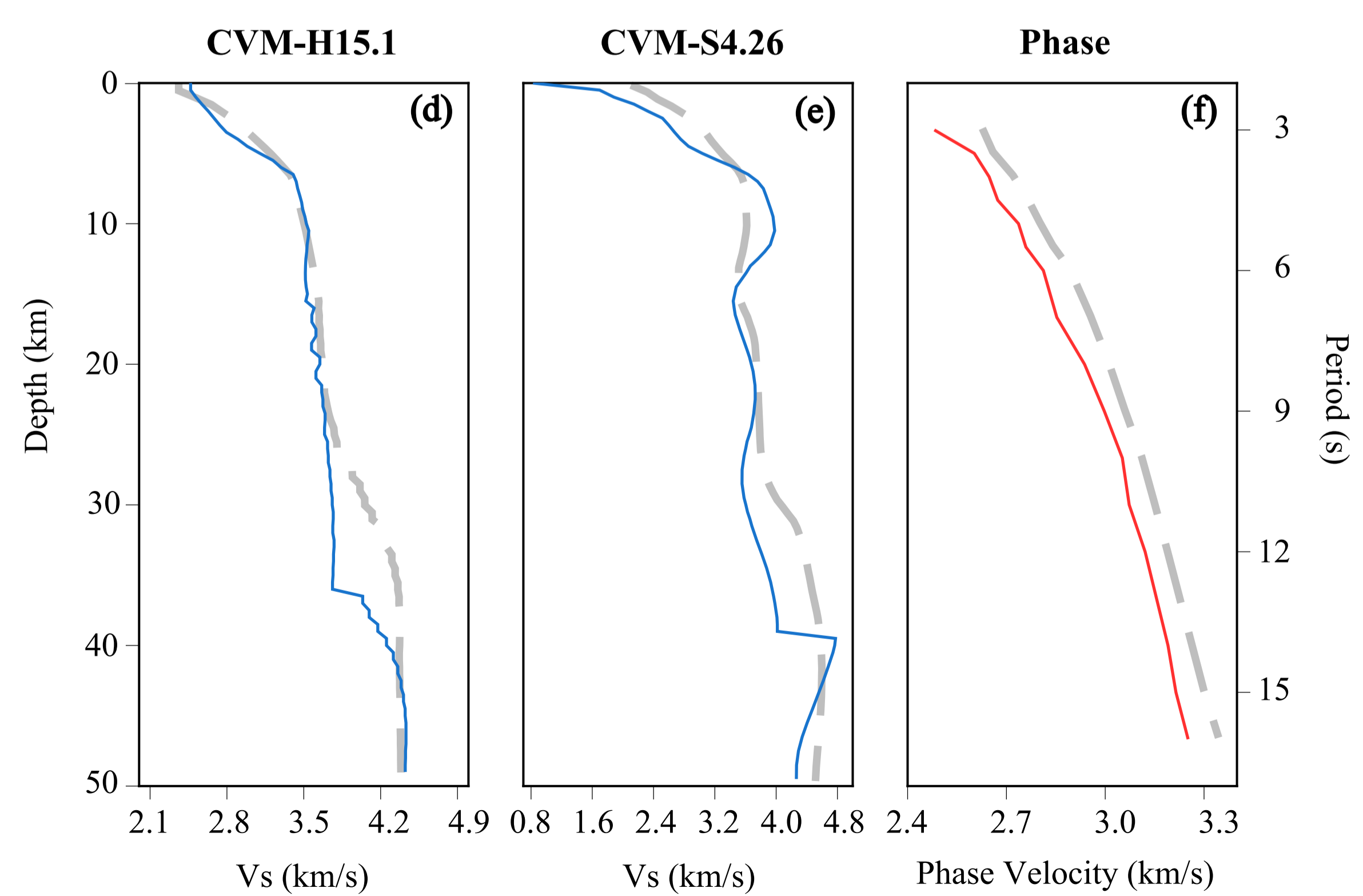
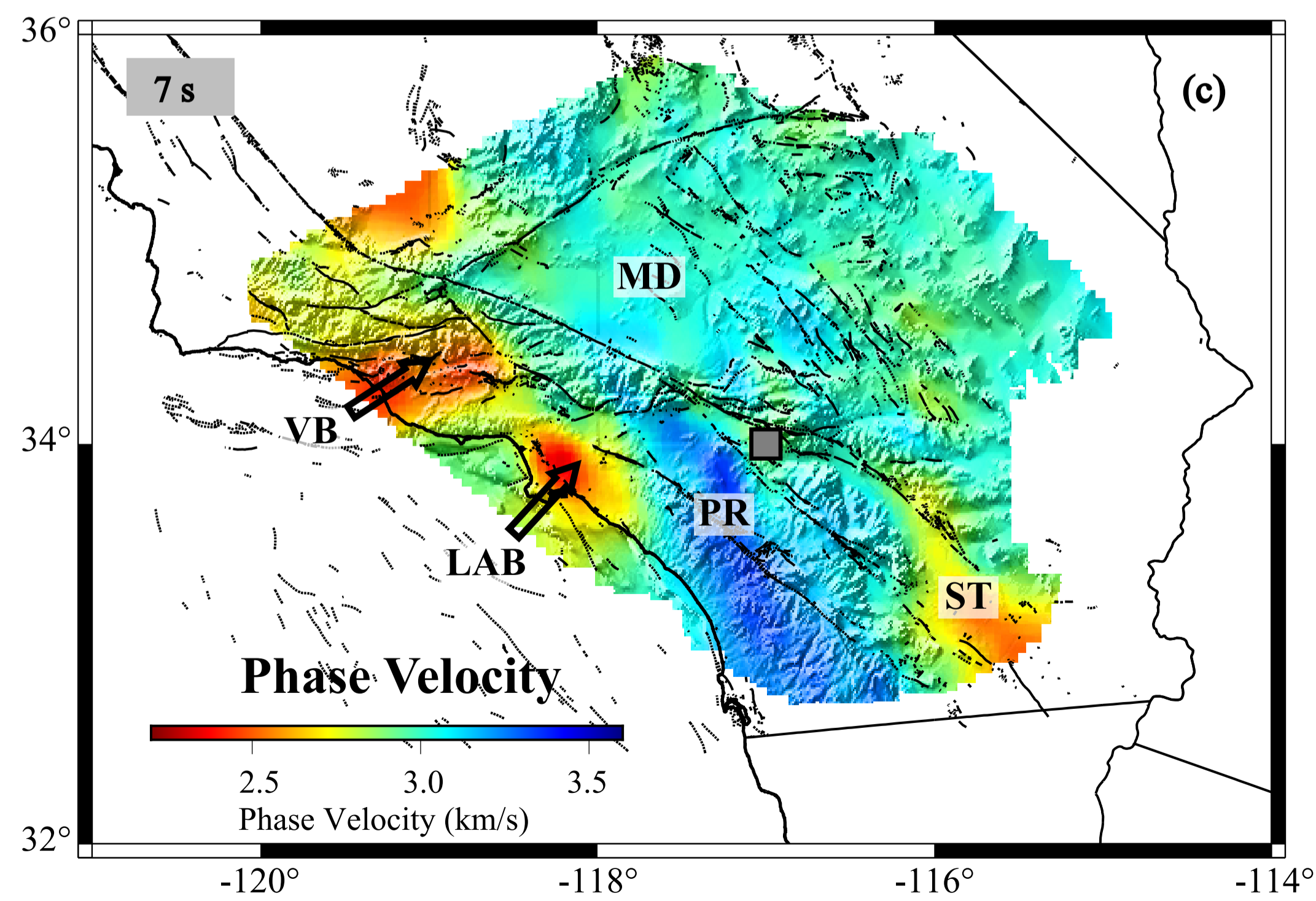
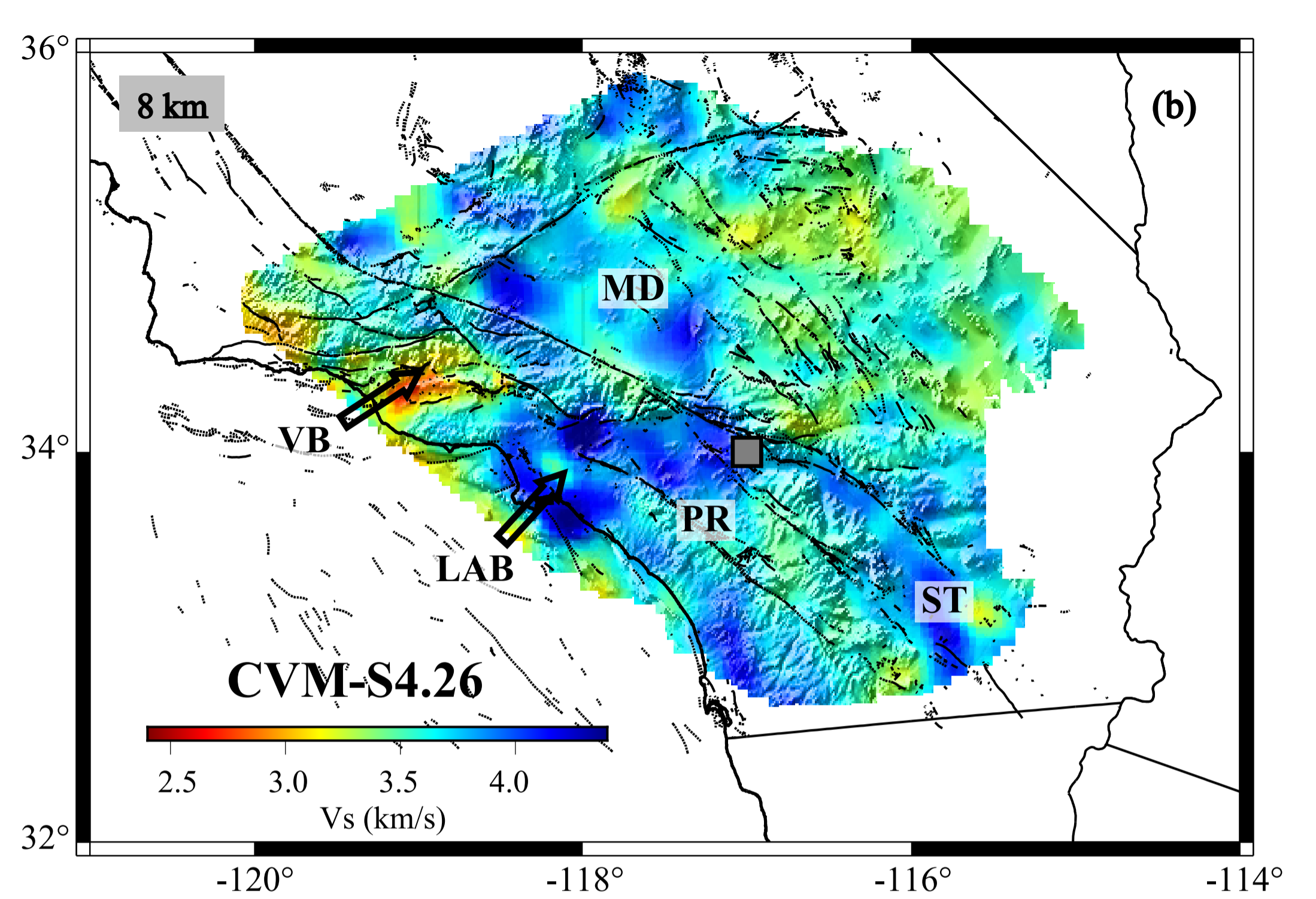
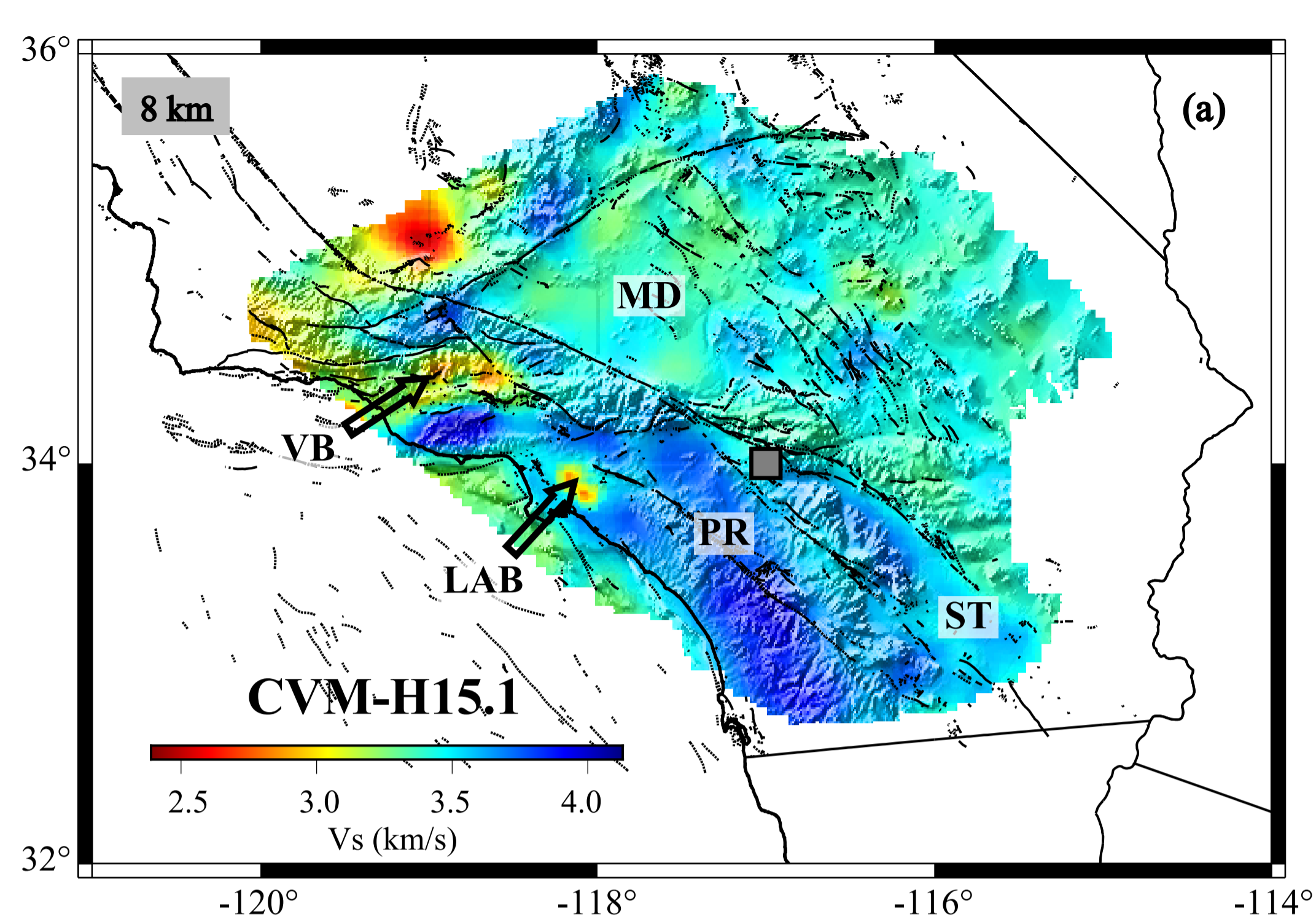


Figure 4.

