Electronic Supplementary Material (ESI) for Physical Chemistry Chemical Physics. This journal is © the Owner Societies 2023

Pushing the boundaries of VCD spectroscopy in natural product chemistry

Tom Vermeyen,^{a,b} Andrea N. L. Batista,^c Alessandra L. Valverde,^c Wouter Herrebout *^a and João M. Batista Jr.*^d

 ^aDepartment of Chemistry, University of Antwerp, Groenenborgerlaan 171, B-2020 Antwerp, Belgium. E-mail: <u>wouter.herrebout@uantwerpen.be</u>
 ^bDepartment of Chemistry, Ghent University, Krijgslaan 281, B-9000 Ghent, Belgium.
 ^cInstitute of Chemistry, Fluminense Federal University, Outeiro de São João Batista s/n, 24020-141 Niterói-RJ, Brazil.
 ^dFederal University of São Paulo, Institute of Science and Technology, R. Talim 330, 12231-280, São José dos Campos-SP, Brazil. E-mail: <u>batista.junior@unifesp.br</u>

Electronic Supplementary Information

Summary

Experimental detais	2
Results of IR/VCD visual inspection	2
Figure S1. IR and VCD spectra of pinane type monoteperpenes	6
Figure S2. IR and VCD spectra of menthane type 1 monoteperpenes	7
Figure S3. IR and VCD spectra of menthane type 2 monoteperpenes	8
Figure S4. IR and VCD spectra of bornane type monoteperpenes	9
Figure S5. IR and VCD spectra of fenchane type monoteperpenes	10
Figure S6. IR and VCD spectra of geraniol type monoteperpenes	11
Figure S7. IR and VCD spectra of other types	12
Figure S8. IR and VCD spectra of the mixture of pinane type monoteperpenes	13
Figure S9. IR and VCD spectra of the mixture of menthane type 1 monoteperpenes	13
Figure S10. IR and VCD spectra of the mixture of menthane type 2 monoteperpenes	14
Figure S11. IR and VCD spectra of the mixture of bornane type monoteperpenes	14
Figure S12. IR and VCD spectra of the mixture of fenchane type monoteperpenes	15
Figure S13. IR and VCD spectra of the mixture of geraniol type monoteperpenes	15
Figure S14. GC-MS analysis of tea tree oil	16
Figure S15. GC-MS analysis of rosemary oil	17
Figure S16. GC-MS analysis of lavender oil	18

Figure S17. GC-MS analysis of ylang-ylang oil	19
Machine Learning model structure and development	20
Figure S18. Similarity of IR spectra for each pair of monoterpenes	21
Figure S19. Similarity of VCD spectra for each pair of monoterpenes	21
Technical details of Machine Learning model and hyperparameter optimisation	22
Figure S20. Optimization of regularization strength for VCD and IR	22
Predictions by Machine Learning model on mixtures of known composition	23
Table S1. Content of the experimental mixtures	24
Figure S21. Predicted concentrations for the chiral terpenes by the VCD-based model on the commixtures A-F with pinane, menthane 1 and menthane 2 types	bination of 26
Figure S22. Predicted concentrations for the chiral terpenes by the VCD-based model on the commixtures A-F with bornane and fenchane types	bination of 27
Figure S23. Predicted concentrations for the chiral/achiral terpenes by the IR-based model on the co of mixtures A-F with pinane, menthane 1 and menthane 2 types	mbination 28
Figure S24. Predicted concentrations for the chiral/achiral terpenes by the IR-based model on the co of mixtures A-F with bornane and fenchane types.	mbination 29
Figure S25. Predicted concentrations for the chiral terpenes on the combination of mixture A-F with r I and J by the VCD-based model	nixtures H, 30
Figure S26. Predicted concentrations for the chiral/achiral terpenes on the combination of mixtur mixtures H, I and J by the IR-based model	e A-F with 31
Figure S27. Coefficients for the L2-regularized VCD model	32
Figure S28. Coefficients for the L2-regularized IR model	33
Predictions by Machine Learning model on oils	34
Table S2. Accuracy of the predictions on the essential oils by the VCD model	35
Table S3. Accuracy of the predictions on the essential oils by the IR model	35

Experimental details

All monoterpenes and essential oils used were purchased from Sigma-Aldrich and used without further purification. The artificial mixtures were prepared by mixing equal amounts of each compound. IR and VCD spectra were recorded simultaneously with a BioTools Chiral/*R*-2x FT-VCD spectrometer with either single or dual-PEM setups using a resolution of 4 cm⁻¹ and a collection time of 10-12 h. The optimum retardation of the ZnSe photoelastic modulator(s) (PEM) was(ere) set at 1400 cm⁻¹. The IR and VCD spectra were recorded in CDCl₃ solutions (0.2-0.8M) in a BaF₂ cell with a 100 µm path length. Minor instrumental baseline offsets were eliminated from the final VCD spectrum by subtracting the VCD spectrum each compound from that obtained for the solvent under the same conditions. The database of VCD and IR spectra is publicly available and can be retrieved using the following DOI [10.5281/zenodo.7875469]. The absolute configuration of each monoterpene when applicable was secured by DFT calculations at the B3PW91/PCM(CHCl₃)/6-311G(d,p) level (data not shown). These calculations also allowed the assignment of the vibrational origin of specific bands.

Results of IR/VCD visual inspection

Figures S1-S7 present the superposition of the IR/VCD spectra of individual monoterpenes within a given molecule type, namely, pinane, menthane 1 and 2, bornane, fenchane and geraniol type as well as the spectra of the single representatives of carene and thujane types along with cineole. Then, IR/VCD spectra of the mixtures of compounds of each type are presented in figures S8-S13. The following discussion about spectral markers are focused on transitions able to tell apart compounds within the same molecule type. Regarding pinane type monoterpenes, the main discriminatory IR bands observed (Fig. S8) were those at 1639 cm⁻¹ present in (S)-(–)- β -pinene (exocyclic double bond stretching); 1616 cm⁻¹ present in (*R*)-(–)-myrtenal and (1*S*)-(–)-verbenone (α , β -unsaturated double bond stretching); 1250 cm⁻¹ present in (1*R*)-(–)-myrtenol and (1*R*,2*R*,3*S*,5*R*)-(–)-pinanediol (C-O stretching); and 1035 cm⁻¹ present in (1*R*,2*R*,3*R*,5*S*)-(–)-isopinocampheol (C-O stretching coupled to C-H bendings of the whole molecular framework). The VCD marker bands included those at (–)-1195 cm⁻¹ present in (S)-(–)- β -pinene (C-H bendings of the whole molecular framework); (+)-1126 cm⁻¹ present in (R)-(+)- α -pinene (C-H bendings of the whole molecular framework); (+)-1035 cm⁻¹ present in (1R,2R,3R,5S)-(-)-isopinocampheol (C-O stretching coupled to C-H bendings of the whole molecular framework); and (-)-967 cm⁻¹ present in (1R)-(-)-myrtenal (C-sp³-C-sp² stretching coupled to C-H bendings of the whole molecular framework). For menthane type 1 molecules (Fig. S9), the IR discriminative bands were those at 1643 cm⁻¹ present in (S)-(-)perillaldehyde, (S)-(-)-perillyl alcohol, (R)-(-)-carvone (broader shoulder) and (R)-(+)limonene (stretching terminal double bond); 1415 cm⁻¹ present in (S)-(–)-perillaldehyde (CH₂ scissoring), 1045 cm⁻¹ present in (S)-(–)- α -terpineol (C-sp³-C-sp² stretching coupled to C-H bendings of the whole molecular framework); and 975 cm⁻¹ present in (S)-(–)-perillyl alcohol (Coupled C-C stretchings and C-H bending of the whole molecular framework). As for VCD marker bands, the band at (-)-1434 cm⁻¹ (asymmetric CH₃ bending and C-H₂ scissoring modes) was present in all molecules, except (R)-(-)-terpinen-4-ol, while that at (-)-1250 cm⁻¹ (C-H bendings of the whole molecular framework and C-H₂ twisting modes) was present in all molecules, except (S)-(–)-perillaldehyde. The band (–)-1045 cm⁻¹ was present only in (S)-(–)- α -terpineol (C-sp³-C-sp² stretching coupled to C-H bendings of the whole molecular framework). For the menthane type 2 monoterpenes (Fig. S10) important IR bands include those at 1677, 1614 (broad) (α , β -unsaturated carbonyl stretchings), and 1286 cm⁻¹ (C-sp²-C sp^2 stretch) present in (R)-(+)-pulegone; 1642 (terminal double bond stretch), 1394 (double bond scissoring), and 1286 cm⁻¹ (coupled O-H and C-H bending modes) present in (1R,2S,5R)-(-)-isopulegol. As for VCD, at around 1286 cm⁻¹, both (R)-(+)-pulegone and (1R, 2S, 5R)-(-)isopulegol presented a positive band, however, in contrast to IR, these bands were better resolved due to their different vibrational origins. At 1103 cm⁻¹ (C-CH₃ and C-O stretchings) a positive VCD band was characteristic of (1R, 2S, 5R)-(-)-menthol, while a +, - band (low to high wavenumbers) centered at 1070 cm⁻¹ (same C-C strecthings coupled to bendings of the whole molecular framework) was observed for (1S, 2R, 5R)-(+)-isomenthol. A negative 1012 cm⁻¹ band (C-C stretchings and C-H isopropyl bending) was observed for both (1R,2S,5R)-(-)isopulegol and (1S,2S,5R)-(+)-neomenthol, while a positive band at 962 cm-1 was present in the spectra of (1R,2S,5R)-(-)-isopulegol, (1S,2R,5R)-(+)-isomenthol, and (1S,2S,5R)-(+)neomenthol. While the band at (-)-1012 cm⁻¹ band seems to be selective of menthane molecules with trans relationship between the isopropyl and methyl groups, the (+)-962 cm⁻ ¹ band arise from C-C stretchings and C-H bendings of the whole molecular framework, being representative of the menthane type 2 scaffold. Considering bornane type molecules (Fig. S11) the IR marker bands identified include that at 1415 cm⁻¹ observed for (1*R*)-(+)-camphor (CH₂ scissoring in the vicinity of carbonyl group); those at 998 and 1068 cm⁻¹ present in (\pm) isoborneol (C-C-O stretching coupled to CH₂ rocking vibrations), and those at 1012, 1229 and 1253 characteristic of (S)-(-)-endo-borneol (C-C-O stretching coupled to CH₂ rocking vibrations). These latter vibrations reflect the endo and exo orientations of the OH group in these stereoisomers. Distinctive VCD bands in bornane type molecules were observed at (+)-1320 and (+)-1166 cm⁻¹ for (1*R*)-(+)-camphor (C-C stretch of quaternary bridgehead carbon coupled to CH₂ wagging and C-C stretch of quaternary bridge carbon coupled to methyne bending, respectively); at 1259 cm⁻¹ a negative couplet-like band (from low to high wavenumbers) was observed for (S)-(-)-endo-borny acetate (C-sp²-O stretching coupled to CH₂ wagging and methyne bending modes); centered at 1125 cm⁻¹ a negative couplet-like band (from low to high wavenumbers) was observed for (S)-(-)-endo-borneol and (S)-(-)endo-borny acetate (C-O stretching coupled to C-C stretches of the whole molecular framework and methyne bending modes); at 1070 and 981 cm⁻¹ two positive VCD bands were observed for (S)-(-)-endo-borny acetate (C-C stretching coupled to C-H bendings involving the whole molecular framework), while positive VCD bands at 1053 and 981 cm⁻¹ were present

for (S)-(–)-endo-borneol. Interestingly, the bands at 1135, 1070 and 981 cm⁻¹ (fundamentals 124, 116, and 101, respectively in the original publication) could have been used to assign the absolute configuration of the monoterpenic portion of the monoterpene chromane esters isolated from Peperomia obtusifolia in 2011 (J. Org. Chem. 2011, 76, 2603-2612). At that time, the stereochemistry of the bornyl moieties tethered to the 3,4-dihydro-5-hydroxy-2,7dimeth-yl-8-(3"-methyl-2"-butenyl)-2-(4'-methyl-1',3'-pentadienyl)-2H-1-benzopyran-6carboxylyc acid were determined using arithmetic operations on experimental and calculated spectra for diastereomeric compounds. In the case of fenchane type molecules (Fig. S12), the IR bands at 1080, 1064 and 1010 cm⁻¹ (C-C stretchings coupled to C-H bendings involving the whole molecular framework) were present in (1R)-(+)-endo-fenchyl alcohol, while the band at 1023 cm⁻¹ (C-sp³-C-sp² stretching coupled to C-H bendings involving the whole molecular framework) was characteristic of (S)-(+)-fenchone in this region. In the VCD spectra, the positive band at 1080 cm⁻¹ was observed for (1*R*)-(+)-*endo*-fenchyl alcohol, while the (+)-1023 cm⁻¹ and (–)-996 cm⁻¹ (C-H bendings involving the whole molecular framework), were characteristic of (S)-(+)-fenchone. Finally, considering the linear terpenes (geraniol type) (Fig. S13), the IR band at 1672 cm⁻¹ was observed for all molecules since it involved the stretching of the trisubstituted double bound from the terminal isoprene unit. Bands at 1637 and 1412 cm⁻¹ were observed for (*R*)-(–)-linalyl acetate, (*R*)-(–)-linalool and (*S*)-(+)- β -citronellene and involved stretching and scissoring modes of their terminal double bond; at 1477 cm⁻¹ a shoulder band was present only in the spectrum of $(S)-(-)-\beta$ -citronellol (scissoring of CH₂-OH); at 1106 cm⁻¹ a band was observed for (R)-(–)-linalool arising C-O stretching and O-H bending of the tertiary alcohol, while the same vibration modes were observed at 1054 cm⁻¹ for the primary alcohol (S)-(–)- β -citronellol. Despite the lower intensities and noisier VCD spectra observed for linear monoterpenes, some discriminatory VCD bands were identified, such as the (+)-1089 cm⁻¹ observed for (S)-(+)- β -citronellene (C-H bendings involving the whole molecular framework); the (–)-1075 cm⁻¹ band (C-C stretches coupled to C-H and O-H bending modes) observed for (R)-(-)-linalool, and the (+)-1054 cm⁻¹ observed for (S)-(-)- β citronellol.



Figure S1. Superposition of IR/VCD experimental spectra in CDCl₃ of pinane type monoterpenes. (Black) (1*R*)-(–)-myrtenol; (Green) (1*R*)-(–)-myrtenal; (Red) (1*R*)-(–)-myrtenyl acetate; (Blue) (*S*)-(–)- β -pinene; (Cyan) (*R*)-(+)- α -pinene; (Magenta) (1*R*,2*R*,3*S*,5*R*)-(–)-pinanediol; (Yellow) (1*S*)-(–)-verbenone; (Navy) (1*S*,2*S*,5*S*)-(–)-2-hydroxy-3-pinanone; (Dark Yellow) (1*R*,2*R*,3*R*,5*S*)-(–)-isopinocampheol.



Figure S2. Superposition of IR/VCD experimental spectra in $CDCI_3$ of menthane type 1 monoterpenes. (Black) (*S*)-(–)- α -terpineol; (Green) (*S*)-(–)-perillyl alcohol; (Red) (*R*)-(–)-terpinen-4-ol; (Blue) (*R*)-(+)-limonene; (Cyan) (*S*)-(–)-perillaldehyde; (Magenta) (*R*)-(–)-carvone.



Figure S3. Superposition of IR/VCD experimental spectra in CDCl₃ of menthane type 2 monoterpenes. (Black) (1R,2S,5R)-(-)-menthol; (Green) (1R,2S,5R)-(-)-isopulegol; (Red) (1S,2S,5R)-(+)-neomenthol; (Blue) (1S,2R,5R)-(+)-isomenthol; (Cyan) (R)-(+)-pulegone.



Figure S4. Superposition of IR/VCD experimental spectra in $CDCl_3$ of bornane type monoterpenes. (Black) (1*R*)-(+)-camphor; (Green) (*S*)-(–)-*endo*-borneol; (red) (*S*)-(–)-*endo*-borny acetate; (Blue) (±)-isobornyl acetate (IR only); (Cyan) (±)-isoborneol (IR only). Gap in the carbonyl region due to high noise level.



Figure S5. Superposition of IR/VCD experimental spectra in $CDCl_3$ of fenchane type monoterpenes. (Black) (*S*)-(+)-fenchone; (Red) 1R)-(+)-*endo*-fenchyl alcohol.



Figure S6. Superposition of IR/VCD experimental spectra in $CDCI_3$ of geraniol type monoterpenes. (Black) (*S*)-(–)- β -citronellol; (Green) (*R*)-(–)-linalool; (Red) (*R*)-(–)-linalyl acetate; (Blue) (*S*)-(+)- β -citronellene; (Cyan) (*S*)-(–)- β -citronellal. Gap in the carbonyl region due to high noise level.



Figure S7. Superposition of IR/VCD experimental spectra in CDCl₃ of: (Black) (1*S*)-(+)-3-carene; (Red) ((1*S*,4*R*)-(-)- α -thujone; (Blue) cineole (IR only). Gap in the carbonyl region due to high noise level.



Figure S8. Monoterpenes identified from the artificial mixture of pinane type molecules by means of visual IR/VCD spectral markers. Selected vibrational frequencies and molecular origin also provided.



Figure S9. Monoterpenes identified from the artificial mixture of menthane type 1 molecules by means of visual IR/VCD spectral markers. Selected vibrational frequencies and molecular origin also provided. Shaded areas represent common bands.



Figure S10. Monoterpenes identified from the artificial mixture of menthane type 2 molecules by means of visual IR/VCD spectral markers. Selected vibrational frequencies and molecular origin also provided.



Figure S11. Monoterpenes identified from the artificial mixture of bornane type molecules by means of visual IR/VCD spectral markers. Selected vibrational frequencies and molecular origin also provided. Shaded areas indicate couplet signals.



Figure S12. Monoterpenes identified from the artificial mixture of fenchane type molecules by means of visual IR/VCD spectral markers. Selected vibrational frequencies and molecular origin also provided.



Figure S13. Monoterpenes identified from the artificial mixture of geraniol type molecules by means of visual IR/VCD spectral markers. Selected vibrational frequencies and molecular origin also provided. Shaded areas represent common bands.



Figure S14. GC-MS analysis of tea tree oil.



Figure S15. GC-MS analysis of rosemary oil.

			Childh	natogram Laven	der on ee	Jennessonation De	add rojeerro		,u			TIC
4,395,999	3.230	0.140										
	-	0										
1												
1												
717	500	F	4									
89 27 880	994	1.65	630 26.9 549									
L I I I I	訂正顧	<u> </u>	58 128									
10.0	0	20.0	30.0	40.0		50.0	60.0	7	0.0	80.0	90.0	97.
												min
					Pe	ak Report						
Peak#	R.Time	I.Time	F.Time	Area	Area%	Height	Height%	A/H	Mark	Name		1
1	0.080	0.035	6./35	289075	0.83	124599	1.24	2.32		(1S)-2,6,6-1rimethy	1bicyclo[3.1.1]	nept-2-e
3	8 717	8.080	8 815	731104	2 11	274715	2 72	2.59		beta - Myrcene		
4	9.399	9.360	9.445	81045	0.23	30661	0.30	2.64		(1S)-2.6.6-Trimethy	lbicvclo[3.1.1]	hept-2-e
5	9.964	9.925	10.025	99883	0.29	39574	0.39	2.52		Benzene, tert-butyl-	5 1. 1	1
6	10.138	10.100	10.155	59888	0.17	26581	0.26	2.25		Cyclopropane, 1,1'-	ethenylidenebis	-
7	10.197	10.155	10.270	189289	0.55	60663	0.60	3.12	V	5-Hepten-2-one, 6-r	nethyl-	
8	10.577	10.545	10.645	102684	0.30	39530	0.39	2.60		(1S)-2,6,6-Trimethy	lbicyclo[3.1.1]	hept-2-e
10	13 230	13 100	13 305	15465248	44.65	4246477	42.12	3.02		Linalool	-dimensi-	
10	15.927	15.870	16.000	394582	1.14	122968	1.22	3.21		Isoborneol		
12	16.135	16.085	16.190	44451	0.13	15860	0.16	2.80		3-Methyl-3-nitrobut	-1-ene	
13	16.473	16.410	16.550	340825	0.98	106100	1.05	3.21		Terpinen-4-ol		
14	17.106	17.070	17.165	43573	0.13	16647	0.17	2.62		4-Pentene-2-ol, 2-m	ethyl	
15	20.140	20.020	20.225	14775784	42.66	4393637	43.58	3.36		Linalyl acetate		
16	21.657	21.595	21.750	453468	1.31	124372	1.23	3.65		.betaMyrcene	7 dimentional	
1/	24.813	24./50	24.925	3/39/2	1.09	99104 50077	0.98	3.79		2,0-Octadien-1-ol, 2	., / -dimethyl-	
10	25.050	25.585	23.723	546320	1.58	154090	1.59	3.03		(Z Z)- alpha -Farnes	ene	
20	28.649	28.600	28.725	107149	0.31	32973	0.33	3.25		(E)-,beta,-Famesene		
20		20:000	201720	34639731	100.00	10082502	100.00	5.25		() ioetai i antesene		

Figure S16. GC-MS analysis of lavender oil.



Figure S17. GC-MS analysis of ylang-ylang oil.

Machine Learning model structure and development

Due to the absence of a large monoterpene mixture dataset, a set of in-silico mixture IR and VCD spectra was generated. These spectra were constructed as random linear combinations of the monoterpene spectra, upon which gaussian noise is added. As the in-silico mixture spectra are linear combinations, a (L2-regularised) linear model was chosen as the basis for the ML model. The model was trained on the VCD and IR in-silico mixture spectra separately to predict the concentration of each monoterpene. During training, the model teaches itself the marker bands for each terpene, identifies which compounds can attenuate their intensity and from which areas of the spectrum non-marker bands can improve detection. The added noise guided the model to ignore spectral features with intensities close to the noise level and the regularisation implored the model to focus on the wavenumber most important for detecting the specified terpene. By doing so, we limited the overfitting of the model to the in-silico spectra. The noise level was based on the noise level found in the experimental IR and VCD spectra. The strength of regularisation was increased as much as possible without significantly decreasing the accuracy of the in-silico concentration predictions (R^2 of approximately 0.98 for unseen in-silico spectra). Technical details on the training and optimization procedure are provided in the next section.

Prior to evaluating the obtained results on the experimental mixtures, we discuss the differences in diversity of the IR and VCD spectra for the monoterpenes. As shown in Figure S18, the IR spectra are less diverse and grouped into 3 clusters: compounds containing a nonconjugated carbonyl group, a conjugated carbonyl group or lacking any carbonyl group moiety. The spectra within each cluster are strongly similar, increasing the difficulty in separating the contributions of individual terpenes. The low noise level of IR can compensate for the lower diversity, as small contributions can be more easily discerned. In contrast to IR, the VCD spectra are much less correlated as shown in Figure S19. The individual contributions of different chiral terpenes are, therefore, expected to be more easily separated from each other. The higher noise level of VCD and baseline uncertainties could increase the difficulty of detecting all contributions, though. The VCD-based model holds two additional advantages for analysis of complex mixtures. The transparency of VCD to achiral compounds improves the stability of the model towards the presence of achiral compounds absent from the dataset. Also, the high sensitivity of VCD to molecular chirality introduces said sensitivity in the model, enabling future use of the model for determination of stereochemistry of essential oil components.



Figure S18. Similarity of IR spectra for each pair of monoterpenes. Similarity is expressed as the absolute cosine similarity. The order of the monoterpenes is based on hierarchical clustering on the IR similarity values.



Figure S19. Similarity of VCD spectra for each pair of chiral monoterpenes. Similarity is expressed as the absolute cosine similarity. The order of the monoterpenes is based on hierarchical clustering on the VCD similarity values.

Technical details of Machine Learning model and hyperparameter optimisation

As mentioned in the previous section, the basis of the ML model is a L2-regularized linear regression (also known as Ridge regression) and the model is trained to predict the concentrations of each terpene from the noisy in-silico mixture spectra. The output of the VCD model is a 33-dimensional vector containing the concentration of each chiral terpene. For the IR model the output is a 36-dimensional vector containing the concentrations of all chiral and achiral terpenes. The spectral intensities (IR or VCD) for each of the 441 wavenumbers between 950 and 1800 cm⁻¹ constitute the input of the model. The model was built and trained using the scikit-learn library (version 0.24.2) [L. Buitinck, G. Louppe, M. Blondel, F. Pedregosa, A. Mueller, O. Grisel, V. Niculae, P. Prettenhofer, A. Gramfort, J. Grobler, R. Layton, J. Vanderplas, A. Joly, B. Holt and G. Varoquaux, 2013, arXiv:cs.LG] and default settings were used unless specified otherwise.

For the model, the strength of the regularization, referred to as α , is an important hyperparameter requiring optimization. Both the VCD and IR model were trained with a range of α values using 10-fold cross validation. The resulting performance for the in-silico training and validation sets are shown in Fig S20. For smaller α values the VCD model is overfitted to the training set and larger α values result in underfitting. By setting α to 1.10⁻⁹, both influences are balanced and the resulting model predicts the terpene concentrations with a R² of ± 0.98 for in-silico mixtures. For the IR model, we chose the largest α value (1.10⁻¹) that resulted in a similar accuracy (R² of ± 0.98). By doing so, we keep the relative level of regularisation consistent for the VCD and IR models.



Fig S20: Optimization of regularization strength α for the VCD (left panel) and IR (right panel) in-silico mixture spectra. The reported R² values are the averages of R² for each cross-validation fold and the error on this average is the standard deviation for the R² values.

After the hyperparameter optimization of the VCD model, the VCD models arising from each fold are combined into an ensemble where the predicted concentration for a single terpene is the mean value of the predicted concentrations of each model and the standard deviation is used to quantify the error upon the mean value. This approach is known as bagging [L. Breiman, *Mach. Learn.* 1996, **24**, 123–140.] and can improve the robustness of the predictions while providing a notion for the uncertainty upon the predicted values. The approach is repeated for the IR model using the IR models of each fold.

Predictions by Machine Learning model on mixtures of known composition

The contents of the different experimental mixtures are provided in table S1 and the predicted relative concentrations for the IR and VCD models are shown in Figure 3 and Figures S21-S26. Performance of the L2-regularized model on the VCD spectra of mixtures A-F is very promising. For most terpenes, the presence of a specific terpene in a mixture was linked with a higher predicted concentration for that terpene. The VCD-based model could not properly detect the presence/absence of (S)-(-)-perillyl alcohol, (1S)-(+)-carene, (1S,2R,5R)-(+)isomenthol and camphor in mixtures A-E (see Figure 3). Camphor was the only terpene for which both enantiomers are present in a mixture: (1S)-(-)-camphor in mixture A and (1R)-(+)camphor in mixture C. The camphor concentration was expressed in terms of (1R)-(+)camphor for the VCD model so the strong negative prediction should indicate the presence of (1S)-(-)-camphor, as enantiomers have mirror image VCD spectra. So the large negative concentration predicted for mixture A shows that the model has identified (1S)-(-)-camphor. However, no clear detection of (1R)-(+)-camphor was obtained for mixture C. For (R)-(-)linalyl acetate, the largest predicted concentration out of the mixtures corresponds to mixture E. For two other mixtures void of (R)-(-)-linalyl acetate, however, rather large concentrations were predicted. This is likely a consequence of its low VCD intensity. Detecting such low contributions in a mixture spectrum will require large coefficients and the prediction quality will be more easily affected by noise. If training is performed with L1-regularization instead of L2, invoking sparsity in the model, the main difference on model performance lies in that the largest predicted (1S)-(+)-3-carene concentration is obtained for mixture D, while the presence of trans-caryophyllene cannot be detected reliably. While the largest predicted concentrations for each terpene correctly reflects its presence in a mixture, the gap between predicted concentrations when the terpene is present or absent was small for (S)-(-)citronellal, (S)-(+)-β-citronellene, (1S)-(+)-3-carene (for L1) and (R)-(-)-linalyl acetate. The VCD patterns arising from the carbonyl vibration are particularly sensitive to the molecular environment. In complex mixtures, a mixture spectrum could therefore deviate from the linear approximation for the mixture spectra. We trained the linear model again while omitting signals above 1500 cm⁻¹, but performance did not improve.

The IR spectra contain less noise and intensities cannot partially cancel each other, but they are more strongly correlated. The balance between these differences determine the performance of an IR-based model. We trained a L2-regularized model trained on in-silico IR mixture spectra and assessed its performance on the IR spectra of mixtures A-F. The model was trained and validated on detecting the presence of chiral and achiral/racemic terpenes (i.e. cineole, isoborneol and isobornyl acetate). The model could not detect the presence of two terpenes: isomenthol and trans-caryophyllene (sesquiterpene). Also the gap between predicted concentrations for when a terpene is present or absent was small for α -pinene, 3-carene and β -citronellene (see figure 3). The linear approach suggested in this work worked slightly better for IR than for VCD. A combination of the higher noise level, higher uncertainty on the baseline or the possibility of cancelling intensities is likely the reason for this.

The question now remained whether the linear model for a single terpene used only the marker bands or also leveraged the other regions in the spectra to improve its predictions. The coefficients of the L2 linear model are plotted for each terpene in figures S27 and S28. To address this question, we investigated the coefficients from the linear model for a few selected terpenes. For VCD, the model clearly used marker bands to detect some terpenes: e.g. the positive band at 1290 cm⁻¹ for (*R*)-(+)-pulegone, the positive band at 1149 cm⁻¹ for

(S)- $(-)-\alpha$ -terpineol, the positive band at 1052 cm⁻¹ for (S)-(-)-endo-borneol, the negative band at 1718 cm⁻¹ for (S)-(–)-citronellal, and the negative band at 1738 cm⁻¹ for (1R)-(+)-camphor were all heavily used by the respective linear models. The linear model used these marker bands but did not completely rely on them; for many terpenes, numerous non-zero coefficients were found to contribute to their detection. In IR, the model weighed the carbonyl region as important for more terpenes compared to VCD. The coefficients of the IRbased model for carvone provided a clear example of how non-marker bands supplement the marker bands for its detection. For carvone, a strong positive and negative coefficient was observed at 1660 and 1620 cm⁻¹, respectively. Carvone has a strong marker IR band at 1660 cm⁻¹, however so do myrtenal and verbenone. The IR spectra of myrtenal and verbenone both contain a smaller IR band at 1620 cm⁻¹, while carvone does not absorb at this frequency. Thus, the model leveraged the IR intensities at 1620 cm⁻¹ that detected the false positives of myrtenal and verbenone for detecting carvone with the 1660 cm⁻¹ marker band. For pulegone, the most intense IR band at 1677 cm⁻¹ was mainly ignored as multiple terpenes absorb at a similar frequency. The 1614 cm⁻¹ band is notably broad, with the 1610-1560 cm⁻¹ section of the band overlapping only partially with fenchone. The model leveraged all intensities between 1610 and 1560 cm⁻¹ to detect pulegone along with the 1210 and relatively isolated 1288 cm⁻¹ bands.

Next, we tested the performance of the models on the mixtures of pinane type, menthane type 1, menthane type 2, bornane type and fenchane type (Figures S21-S24). For each of these mixtures, we bundled its predictions with those for mixtures A-F and observed whether the presence of a specific terpene was still linked with a higher predicted concentration. The terpenes for which mismatches between predicted concentrations and their presence were already obtained with A-F will not be discussed, but will still be highlighted in the figures. For the mixture of pinane derivatives, the presence of (S)-(-)citronellal and (-)-trans-caryophyllene was wrongly predicted with the VCD model. The IR model wrongly detected pinanediol and limonene. For the menthane 1 mixture, both models could not properly detect (R)-(-)-terpinen-4-ol. The IR spectrum of the menthane 2 mixture allowed to identify all terpenes. The VCD spectrum identified all terpenes present, but a mismatch was obtained for (R)-(-)-carvone and (S)-(-)-citronellal. Interestingly, large positive concentrations were predicted on both spectra for (1S, 2R, 5R)-(+)-isomenthol for which mismatches were obtained on A-F. On the bornane type mixture no additional mismatches were noted for IR and a single wrong detection for (R)-(-)-linally acetate in VCD was noted. For the mixture of fenchone and fenchol, a single wrong prediction was obtained for citronellal on the IR spectrum. The VCD-based model wrongly detected the presence of (R)-(-)-linalyl acetate, (S)-(–)-citronellal and (S)-(+)- β -citronellene. For this set of mixtures, each composed of structurally similar terpenes, the accuracy of the ML approach remained similar to the accuracy obtained for A-F, with on average 1 and 2 additional wrong detections for IR and VCD respectively. The models show clear potential for the analysis of terpene mixtures. Now the question remained how far the application area can be pushed. Therefore, we increased the complexity of the mixtures even further and tested whether the models could still identify the terpenes present. The three mixtures of increased complexity (H-J) are composed of 16-22 terpenes each. The same methodology was repeated, combining the A-F predictions along with each of these three mixtures separately, and the obtained results are shown in Figures S25-S26.

The predictions for J remained fairly accurate, with 3 additional mismatches for (S)-(-)endo-bornyl acetate, (S)-(-)-perillaldehyde and (R)-(+)-pulegone on the VCD spectrum and two mismatches for borneol and limonene on the IR spectrum. For mixture H, the accuracy of the model decreased with 7 additional mismatches for the IR and 4 for VCD. Similarly, a lower accuracy was obtained for mixture I with 9 additional mismatches for IR and 8 for VCD. For mixtures of such complexity, where each terpene provides only a tiny contribution to the mixture spectrum, accurately detecting the terpenes present becomes more challenging. Depending on the exact mixture composition, the models can still extract the presence of the monoterpenes as demonstrated with mixture J.

Mixture	Terpenes present
A	$(1S,4R)-(-)-\alpha$ -thujone, $(1R,2R,3R,5S)-(-)$ -isopinocampheol, $(R)-(+)-\alpha$ -pinene, $(1R)-(+)-endo$ -fenchyl alcohol, $(1S,2S,5R)-(+)$ -neomenthol, $(1S)-(-)$ -camphor, $(1R)-(-)$ -myrtenyl acetate (not present in in-silico mixtures)
В	(S)-(-)-β-pinene, (S)-(-)- <i>endo</i> -borneol, (1R,2S,5R)-(-)-isopulegol, (R)-(+)-α-pinene, (S)- (+)-fenchone, (R)-(+)-limonene, (R)-(+)-pulegone
С	(1 <i>S</i> ,2 <i>S</i> ,5 <i>S</i>)-(–)-2-hydroxy-3-pinanone, (1 <i>R</i> ,2 <i>S</i> ,5 <i>R</i>)-(–)-isopulegol, (1 <i>R</i> ,2 <i>S</i> ,5 <i>R</i>)-(–)-menthol, (1 <i>R</i>)-(–)-myrtenal, (<i>S</i>)-(–)-perillyl alcohol, (1 <i>R</i> ,2 <i>R</i> ,3 <i>S</i> ,5 <i>R</i>)-(–)-pinanediol, (1 <i>R</i>)-(+)- camphor, (<i>R</i>)-(–)-carvone
D	Cineole, (S)-(–)- α -terpineol, (1R)-(–)-myrtenol, (S)-(–)-perillaldehyde, (1S)-(–)-verbenone, (1S)-(+)-3-carene, (1S,2R,5R)-(+)-isomenthol
E	(R)-(–)-linalyl-acetate, (S)-(–)-β-citronellol, (S)-(–)-citronellal, (R)-(–)-linalool, (S)-(+)-β- citronellene
F	(1 <i>S</i> ,2 <i>S</i> ,5 <i>S</i>)-(–)-2-hydroxy-3-pinanone, (<i>S</i>)-(–)-citronellal, (1 <i>R</i> ,2 <i>S</i> ,5 <i>R</i>)-(–)-isopulegol, (<i>S</i>)-(–)-perillyl-alcohol, (–)- <i>trans</i> -caryophyllene, (1 <i>S</i>)-(–)-verbenone, (1 <i>S</i>)-(+)-3-carene, (<i>R</i>)-(–)-carvone
н	$(1S,2S,5S)$ - $(-)$ -2-hydroxy-3-pinanone, (S) - $(-)$ - α -terpineol, (S) - $(-)$ - β -citronellol, (S) - $(-)$ - β -pinene, (R) - $(-)$ -linalool, $(1R)$ - $(-)$ -myrtenal, $(1R)$ - $(-)$ -myrtenol, (S) - $(-)$ -perillaldehyde, (S) - $(-)$ -perillyl-alcohol, (R) - $(+)$ - α -pinene, (S) - $(+)$ - β -citronellene, $(1S)$ - $(+)$ -carene, (S) - $(+)$ -fenchone, $(1S,2S,5R)$ - $(+)$ -neomenthol, (R) - $(-)$ -carvone, (R) - $(+)$ -limonene, (R) - $(+)$ -pulegone, cineole
1	(±)-isobornyl-acetate, (R)-(–)-linalyl-acetate, (1 S ,2 S ,5 S)-(–)-2-hydroxy-3-pinanone, (S)-(–)- α -terpineol, (S)-(–)- β -citronellol, (S)-(–)- β -pinene, (S)-(–)- <i>endo</i> -bornyl acetate, (R)-(–)-linalool, (1 R)-(–)-myrtenal, (1 R)-(–)-myrtenol, (S)-(–)-perillaldehyde, (S)-(–)-perillyl-alcohol, (R)-(+)- α -pinene, (S)-(+)- β -citronellene, (1 S)-(+)-3-carene, (S)-(+)-fenchone, (1 S ,2 S ,5 R)-(+)-neomenthol, (R)-(–)-carvone, (R)-(+)-limonene, (R)-(+)-pulegone, (R)-(–)-terpinen-4-ol, cineole
1	Cineole, (S) - $(-)$ - α -terpineol, (S) - $(-)$ - β -pinene, (S) - $(-)$ - <i>endo</i> -borneol, (S) - $(-)$ -bornyl acetate, $(1R,2S,5R)$ - $(-)$ -isopulegol, $(1R)$ - $(-)$ -myrtenal, (S) - $(-)$ -perillaldehyde, (S) - $(-)$ -perillyl alcohol, (R) - $(+)$ - α -pinene, (S) - $(+)$ -fenchone, $(1S,2R,5R)$ - $(+)$ -isomenthol, $(1S,2S,5R)$ - $(+)$ -neomenthol, (R) - $(-)$ -carvone, (R) - $(+)$ -limonene, (R) - $(+)$ -pulegone

Table S1. Content of the experimental mixtures added for evaluation of the linear models.



Figure S21. Predicted concentrations for the chiral terpenes by the VCD-based model on the combination of mixtures A-F with, from top to bottom, pinane type, menthane 1 type, menthane 2 type respectively. The predicted concentration is colored according to whether the terpene is present (green) or absent (red) for a mixture. The predicted concentrations are highlighted for terpenes when no correct decision boundary can be drawn.



Figure S22. Predicted concentrations for the chiral terpenes by the VCD-based model on the combination of mixtures A-F with, from top to bottom, bornane type, fenchane type, respectively. The predicted concentration is colored according to whether the terpene is present (green) or absent (red) for a mixture. The predicted concentrations are highlighted for terpenes when no correct decision boundary can be drawn.



Figure S23. Predicted concentrations for the chiral/achiral terpenes by the IR-based model on the combination of mixtures A-F with, from top to bottom, pinane type, menthane 1 type, menthane 2 type, respectively. The predicted concentration is colored according to whether the terpene is present (green) or absent (red) for a mixture. The predicted concentrations are highlighted for terpenes when no correct decision boundary can be drawn.



Figure S24. Predicted concentrations for the chiral/achiral terpenes by the IR-based model on the combination of mixtures A-F with, from top to bottom, bornane type, fenchane type, respectively. The predicted concentration is colored according to whether the terpene is present (green) or absent (red) for a mixture. The predicted concentrations are highlighted for terpenes when no correct decision boundary can be drawn.



Figure S25. Predicted concentrations for the chiral terpenes on the combination of mixtures A-F with mixtures H (top), I (middle) and J (bottom) by the VCD-based model. The predicted concentration is colored according to whether the terpene is present (green) or absent (red) for a mixture. The predicted concentrations are highlighted for terpenes when no correct decision boundary can be drawn.



Figure S26. Predicted concentrations for the chiral/achiral terpenes on the combination of mixtures A-F with H (top), I (middle) and J (bottom) by the IR-based model. The predicted concentration is colored according to whether the terpene is present (green) or absent (red) for a mixture. The predicted concentrations are highlighted for terpenes when no correct decision boundary can be drawn.



Figure S27. Coefficients for the L2-regularized VCD model for each chiral monoterpene.



Figure S28. Coefficients for the L2-regularized IR model for each monoterpene.

Predictions by Machine Learning model on oils

Both models were then applied for the 4 essential oils using the decision boundaries fine-tuned with mixtures A-F. Predictions for the terpenes for which each model yielded unreliable predictions on the mixtures A-F were omitted. An overview of the true positives, false positives and false negatives is provided in table S2 for the VCD model and table S3 for the IR model. For tea tree oil, the VCD model wrongly detected the presence of 3 terpenes (S)-(-)-citronellal, (1R, 2R, 3S, 5R)-(-)-pinanediol and (R)-(+)- α -pinene and did not detect (S)-(+)terpinen-4-ol or the tiny fraction of limonene and β-pinene. The IR model clearly detected terpinen-4-ol, along with limonene. However, the model detected the absent compounds β citronellol, menthol, α -pinene and 3-carene. In rosemary oil the VCD model correctly identified the presence of (R)-(+)- α -pinene and (S)-(-)- β -pinene present in the oil, whereas the IR model detected α -pinene, β -pinene and α -terpineol present in the oil. Both models remained undecisive concerning the tiny fraction of (1R)-(+)-camphor present. Regarding lavender oil, the VCD model detected the presence of (R)-(-)-linalool, along with the tiny fraction of (R)-(+)- α -pinene, and the IR model identified both linalool and linally acetate. The VCD model also predicted the presence of (S)-(-)-citronellal, (1R,2R,3S,5R)-(-)-pinanediol, (R)-(-)-carvone and (1S,2S,5R)-(+)-neomenthol. The IR model generated a false positive for limonene, and isopinocampheol (and remains indecisive for menthol). Neither model detected the +/-1% of terpinen-4-ol present in the oil. The tiny fraction of β -pinene was not detected by either model. The tiny fraction of α -pinene was barely detected by the VCD model but not by the IR model. The IR model also potentially detected the presence of isoborneol as the predicted concentration exceeds those for mixtures A-F (in which it was absent).

For ylang-ylang oil, (R)-(-)-linalool was detected by both models. False positives were obtained for (S)-(-)- β -citronellal and (R)-(-)-carvone with the VCD model. The IR model wrongly predicted the presence of following compounds: borneol, isopinocampheol, isopulegol, carene, limonene, pulegone. A large concentration was also predicted for isomenthol, for which the presence could not be properly detected for mixtures A-F. While no decision boundaries could be drawn for isoborneol and isobornyl acetate (due to their absence from A-F), large relative concentrations were obtained for them.

Essential oil	True positives	False positives	False negatives
Tea tree oil	/	(S)-(–)-citronellal, (1R,2R,3S,5R)-(–)-pinanediol, (R)-(+)-α-pinene	(S)-(+)-terpinen-4-ol. Tiny fractions: limonene & β-pinene
Rosemary oil	(R)-(+)-α-pinene, (S)-(–)- β-pinene	(R)-(-)-linalyl acetate, (1R,2R,3S,5R)-(-)-pinanediol, (1S)- (+)-3-carene and (R)-(-)-carvone	Tiny fractions: α -terpineol
Lavender oil	(R)-(–)-linalool* Tiny fractions: (R)-(+)-α- pinene	(<i>S</i>)-(–)-β-citronellal, (1 <i>R</i> ,2 <i>R</i> ,3 <i>S</i> ,5 <i>R</i>)- (–)-pinanediol, (<i>R</i>)-(–)-carvone, (1 <i>S</i> ,2 <i>S</i> ,5 <i>R</i>)-(+)-neomenthol	terpinen-4-ol, (S)-(+)-linalyl acetate Tiny fractions: β-pinene
Ylang-ylang oil	(R)-(–)-linalool	(S)-(–)-citronellal and (R)-(–)- carvone	/

Table S2. Accuracy of the predictions on the essential oils by the VCD model.

*Conflicting with visual inspection results

Table S3. Accuracy of the predictions on the essential oils by the IR model.

Essential oil	True positives	False positives	False negatives
Tea tree oil	terpinen-4-ol. Tiny fractions: limonene	β-citronellol, menthol, α-pinene, carene	Tiny fractions: β-pinene
Rosemary oil	α-pinene, cineole, β- pinene.	α -thujone, perillyl alcohol, β -citronellol, camphor.	Isoborneol Tiny fractions: α-terpineol.
Lavender oil	linalool, linalyl acetate.	limonene, isopinocampheol	terpinen-4-ol. Tiny fractions: α -pinene and β -pinene.
Ylang-ylang oil	linalool	borneol, isopinocampheol, isopulegol, carene, limonene, pulegone	/