

Supporting Information for

**Prediction of Oxidant Exposures and Micropollutant
Abatement during Ozonation Using a Machine Learning
Method**

Dongwon Cha[†], Sanghun Park[‡], Min Sik Kim^{†,§}, Taewan Kim[†], Seok Won Hong^{||},

Kyung Hwa Cho^{‡,*}, Changha Lee^{†,*}

[†]*School of Chemical and Biological Engineering, Institute of Chemical Process (ICP), Seoul National University, 1 Gwanak-ro, Gwanak-gu, Seoul 08826, Republic of Korea*

[‡]*School of Urban and Environmental Engineering, Ulsan National Institute of Science and Technology (UNIST), 50 UNIST-gil, Ulju-gun, Ulsan 44919, Republic of Korea*

[§]*Department of Chemical and Environmental Engineering, Yale University, New Haven, CT 06520, United States*

^{||} *Water Cycle Research Center, Korea Institute of Science and Technology (KIST), Hwarangro 14 Gil 5, Seongbuk-gu, Seoul, 02792, Republic of Korea*

*Corresponding authors

Tel.: +82-52-217-2829 (K.H. Cho); +82-2-880-8630 (C. Lee)

Fax: +82-52-217-2819 (K.H. Cho); +82-2-888-7295 (C. Lee)

E-mail: khcho@unist.ac.kr (K.H. Cho); leechangha@snu.ac.kr (C. Lee)

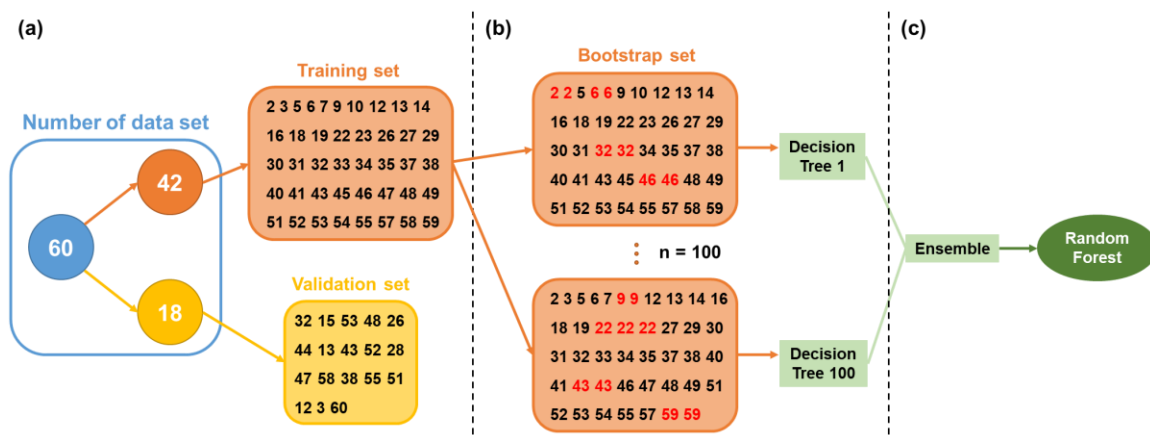
The supporting information contains 28 pages, 9 tables, and 18 figures.

Text S1. Additional details, hyperparameters, and development procedure of the RF model

The RF model was developed using the following hyperparameters after considering the number of data (42 training sets). The minimum parent size was 10 (left as default) and the minimum leaf size was designed as 1 due to the small data set. Additionally, the maximum number of branches (depth) was set to 5 to prevent over-fitting problems driven by the minimum leaf size. Bootstrap aggregating (i.e., bagging) ensemble method was used to improve the stability and accuracy of the model. The number of ensemble learning cycles (i.e., number of trees for ensemble) was set to 100.

The procedure for the development of the RF model is as follows:

- The training and validation data sets were randomly divided into sets of 42 and 18, respectively.
- The training set was resampled with replacements, thus different training data sets were permuted (i.e., bootstrapping). The bootstrapping method can substitute a cross-validation process by producing independent datasets through sampling with replacement. Each bootstrap set was used to train a decision tree model, then 100 tree models were prepared for ensemble modeling.
- Ensemble modeling aggregated the prediction of each tree model and resulted in a final prediction achieving a reduced variance.



Text S2. Comparison of RF and ridge regression models

Two ridge regression models were developed to compare the prediction performance with the RF models (FEEM-Free and FEEM-HighRes) and to evaluate the influence of FEEM features on model performance. The same training and validation data that were used to develop the RF model were also used during the development of the ridge model.

In the case of O₃ exposure prediction, RF-FEEM-Free model shows relatively uniform correlation in training (0.784) and validation (0.732). On the other hand, while Ridge-FEEM-Free model showed a higher correlation in validation (0.836), training fared the worst among FEEM-free models (0.683), with a distinct tendency to underestimate higher O₃ exposure points. For the RF model, the prediction accuracy was increased (from 0.732 to 0.797 in validation) when the FEEM features were included as the input parameters. However, the usage of multiple parameters (i.e., FEEM data) decreased the prediction accuracy (from 0.836 to 0.691 in validation) for the ridge regression model (refer to Figure S14 of the SI)

In the case of •OH exposure prediction, RF-FEEM-Free and RF-FEEM-HighRes models demonstrated a superiority in prediction compared to the ridge regression models regardless of the input parameters (refer to Figure S15 of the SI).

Text S3. Variable importance for FEEM-Free and FEEM-HighRes models

Additionally, variable importance (VI) scores for the FEEM-free model (using only water characteristics) and FEEM-HighRes model were determined. In the case of FEEM-Free model, DOC was the most significant variable for the prediction of O_3 and $\bullet OH$ exposure, while Region 5 of the FEEM contour plot showed the highest importance score when the model included FEEM features (i.e., FEEM-HighRes). As VI scores represent the statistical significance and the contribution of each input parameters with respect to model output, the use of FEEM features and its formidable VI values may result in the increase in prediction accuracy (refer to Figure S17).

Table S1. HPLC analytical conditions for *p*CBA and MPs used in this study

Compound	Flow rate (mL/min)	Eluent composition	UV detection (nm)	Retention time (min)
<i>p</i> CBA	0.8	40 % 0.1 % phosphoric acid, 60 % acetonitrile	234	1.8
Atrazine		50 % 0.1 % phosphoric acid, 50 % acetonitrile	220	1.5
Caffeine		85 % 0.1 % phosphoric acid, 15 % acetonitrile	270	1.2
Carbamazepine		40 % 0.1% phosphoric acid, 60 % acetonitrile	285	1.4
Ibuprofen		40 % 0.1% phosphoric acid, 60 % acetonitrile	215	1.9

Table S2. Oxidant exposures measured for natural waters ($[O_3]_0 = 2.5$ mg/L).

Source	#	Location	O_3 exposure (M s)	$\bullet OH$ exposure (M s)
Nakdong River	1	Uiseong	3.75×10^{-3}	3.48×10^{-10}
	2	Gumi	2.49×10^{-3}	4.69×10^{-10}
	3	Maegok	2.67×10^{-3}	3.19×10^{-10}
	4	Guji	2.63×10^{-3}	3.04×10^{-10}
	5	Namji	2.58×10^{-3}	3.87×10^{-10}
	6	Hanam	1.86×10^{-3}	3.85×10^{-10}
	7	Yangsan	2.29×10^{-3}	3.60×10^{-10}
Han River	8	Chungju	8.91×10^{-3}	5.70×10^{-10}
	9	Yeoju	6.41×10^{-3}	6.80×10^{-10}
	10	Cheongpyeong	1.19×10^{-2}	4.41×10^{-10}
	11	Gwangnaru	7.33×10^{-3}	5.34×10^{-10}
	12	Yangwha	4.83×10^{-3}	4.05×10^{-10}
Geum River	13	Gapcheon	8.51×10^{-3}	4.35×10^{-10}
	14	Daecheong	4.63×10^{-3}	4.52×10^{-10}
	15	Gongju	6.16×10^{-4}	2.46×10^{-10}
	16	Buyeo	7.29×10^{-4}	3.08×10^{-10}
	17	Seocheon	4.04×10^{-3}	3.81×10^{-10}
Yeongsan River	18	Damyang	2.67×10^{-3}	2.87×10^{-10}
	19	Gwangju	3.16×10^{-3}	3.03×10^{-10}
	20	Yeongam	5.10×10^{-3}	2.22×10^{-10}
Seomjin River	21	Sunchang	5.71×10^{-3}	2.00×10^{-10}
	22	Agyang	7.77×10^{-3}	3.03×10^{-10}
	23	Hadong	1.09×10^{-2}	1.36×10^{-10}
Tributaries of the Han River	24	Yangjae	1.55×10^{-2}	5.11×10^{-10}
	25	Tan	1.17×10^{-3}	1.82×10^{-10}
	26	Jungnang	1.92×10^{-3}	2.38×10^{-10}
	27	Cheonggye	1.07×10^{-2}	5.31×10^{-10}
	28	Changreung	8.99×10^{-3}	5.60×10^{-10}
	29	Anyang	5.36×10^{-4}	1.45×10^{-10}
	30	Dorim	1.35×10^{-2}	6.25×10^{-10}

Table S3. Oxidant exposures measured for wastewater effluents ($[\text{O}_3]_0 = 2.5 \text{ mg/L}$).

Region	#	Description	O_3 exposure (M s)	$\bullet\text{OH}$ exposure (M s)
Seoul / Gyeonggi	1	S city, E treatment plant	6.55×10^{-4}	1.96×10^{-10}
	2	S city, T treatment plant	6.15×10^{-4}	1.56×10^{-10}
	3	S city, S treatment plant	4.69×10^{-4}	1.59×10^{-10}
	4	B city, G treatment plant	1.76×10^{-3}	2.45×10^{-10}
	5	B city, Y treatment plant	7.02×10^{-4}	1.49×10^{-10}
	6	P city, P treatment plant	9.25×10^{-4}	9.68×10^{-11}
	7	D city, D treatment plant	1.26×10^{-3}	2.15×10^{-10}
	8	U city, B treatment plant	1.91×10^{-3}	2.60×10^{-10}
South Jeolla	9	Y city	9.12×10^{-3}	4.81×10^{-10}
Gangwon	10	Y county	5.02×10^{-3}	3.44×10^{-10}
	11	Gangwon 1	5.70×10^{-4}	1.31×10^{-10}
	12	Gangwon 2	2.57×10^{-3}	2.51×10^{-10}
	13	Gangwon 3	1.27×10^{-3}	2.49×10^{-10}
	14	Gangwon 4	8.52×10^{-3}	3.47×10^{-10}
	15	Gangwon 5	4.49×10^{-4}	1.13×10^{-10}
	16	Gangwon 6	8.74×10^{-4}	1.36×10^{-10}
	17	Gangwon 7	4.02×10^{-3}	2.87×10^{-10}
	18	Gangwon 8	2.69×10^{-3}	2.18×10^{-10}
	19	Gangwon 9	2.41×10^{-3}	2.09×10^{-10}
	20	Gangwon 10	3.51×10^{-3}	2.91×10^{-10}
	21	Gangwon 11	3.46×10^{-3}	2.93×10^{-10}
	22	Gangwon 12	1.58×10^{-3}	2.58×10^{-10}
	23	Gangwon 13	9.91×10^{-3}	6.40×10^{-10}
	24	Gangwon 14	1.51×10^{-3}	2.11×10^{-10}
	25	Gangwon 15	1.31×10^{-3}	1.60×10^{-10}
	26	Gangwon 16	3.20×10^{-3}	1.79×10^{-10}
	27	Gangwon 17	6.06×10^{-3}	2.19×10^{-10}
	28	Gangwon 18	1.41×10^{-3}	2.14×10^{-10}
	29	Gangwon 19	1.56×10^{-2}	3.51×10^{-10}
	30	Gangwon 20	7.87×10^{-4}	9.38×10^{-11}

Table S4. Water quality parameters (pH, DOC concentration, and alkalinity) for natural waters.

Source	Location	pH	DOC concentration (mg/L)	Alkalinity (mg/L as CaCO ₃)
Nakdong River	Uiseong	8.07	4.58	70
	Gumi	8.94	3.40	70
	Maegok	7.80	3.84	40
	Guji	8.02	4.77	55
	Namji	7.90	3.99	50
	Hanam	8.16	4.59	50
	Yangsan	7.97	4.45	40
Han River	Chungju	8.13	2.21	75
	Yeoju	8.20	2.85	75
	Cheongpyeong	7.62	1.95	35
	Gwangnaru	8.31	2.35	50
	Yangwha	8.11	2.80	50
Geum River	Gapcheon	8.02	2.33	75
	Daecheong	7.87	3.09	35
	Gongju	9.26	3.89	75
	Buyeo	9.08	3.61	75
	Seocheon	8.20	3.56	55
Yeongsan River	Damyang	7.55	4.78	45
	Gwangju	7.72	4.46	50
	Yeongam	8.05	3.73	60
Seomjin River	Sunchang	7.55	4.20	25
	Agyang	7.69	2.68	30
	Hadong	7.66	2.72	30
Tributaries of the Han River	Yangjae	8.09	1.49	70
	Tan	7.89	4.19	90
	Jungnang	7.89	3.40	65
	Cheonggye	8.00	1.66	55
	Changreung	7.81	1.88	50
	Anyang	7.90	4.63	80
	Dorim	7.93	1.74	45

Table S5. Water quality parameters (pH, DOC concentration, and alkalinity) for wastewater effluents.

Region	Description	pH	DOC concentration mg/L	Alkalinity (mg/L as CaCO ₃)
Seoul / Gyeonggi	S city, E treatment plant	7.93	5.27	70
	S city, T treatment plant	7.02	6.01	60
	S city, S treatment plant	7.71	5.48	85
	B city, G treatment plant	7.26	4.29	40
	B city, Y treatment plant	7.17	4.99	50
	P city, P treatment plant	7.97	5.21	210
	D city, D treatment plant	7.58	5.07	70
	U city, B treatment plant	7.51	4.00	40
South Jeolla	Y city	7.76	2.78	40
Gangwon	Y county	7.32	4.02	40
	Gangwon 1	7.62	5.81	75
	Gangwon 2	6.12	3.52	10
	Gangwon 3	7.12	3.79	30
	Gangwon 4	6.73	3.45	10
	Gangwon 5	6.62	4.66	10
	Gangwon 6	7.32	6.06	125
	Gangwon 7	7.20	3.88	35
	Gangwon 8	7.14	6.20	50
	Gangwon 9	6.99	5.07	30
	Gangwon 10	6.80	4.32	15
	Gangwon 11	6.81	4.25	15
	Gangwon 12	7.11	3.95	30
	Gangwon 13	7.20	2.26	20
	Gangwon 14	6.84	4.39	15
	Gangwon 15	7.13	4.83	50
	Gangwon 16	7.25	5.63	50
	Gangwon 17	7.43	3.42	35
	Gangwon 18	7.09	3.67	25
	Gangwon 19	6.24	3.60	5
	Gangwon 20	7.16	5.01	30

Table S6. Percent fluorescence response ($P_{(i,n)}$) values calculated by FRI of the five regions in the FEEM contour plot (for the FEEM-LowRes model) for natural waters.

Source	Location	$P_{(1,n)}$	$P_{(2,n)}$	$P_{(3,n)}$	$P_{(4,n)}$	$P_{(5,n)}$
Nakdong River	Uiseong	1.707	9.152	47.536	7.080	34.526
	Gumi	1.582	10.162	46.710	8.372	33.174
	Maegok	2.005	8.950	47.277	6.774	34.993
	Guji	4.990	19.487	34.973	15.659	24.890
	Namji	4.690	17.555	37.518	13.657	26.581
	Hanam	4.470	17.130	38.349	12.685	27.365
	Yangsan	3.779	13.955	41.833	10.510	29.924
Han River	Chungju	3.301	11.433	44.241	8.478	32.547
	Yeoju	3.412	11.796	44.189	8.874	31.730
	Cheongpyeong	4.672	12.240	44.199	8.902	29.987
	Gwangnaru	4.571	13.862	42.212	9.839	29.516
	Yangwha	4.406	15.641	37.261	13.531	29.161
Geum River	Gapcheon	6.027	15.955	41.475	10.340	26.204
	Daecheong	3.376	10.963	47.460	7.066	31.135
	Gongju	4.253	17.127	38.176	12.630	27.815
	Buyeo	4.692	17.604	38.700	12.413	26.591
	Seocheon	4.430	14.136	43.791	9.262	28.381
Yeongsan River	Damyang	2.589	11.210	43.627	9.063	33.512
	Gwangju	3.315	12.485	42.694	10.034	31.473
	Yeongam	4.128	15.405	39.148	12.675	28.644
Seomjin River	Sunchang	2.087	9.121	44.983	8.422	35.387
	Agyang	3.808	11.844	43.806	9.077	31.465
	Hadong	4.603	12.200	42.811	9.013	31.373
Tributaries of the Han River	Yangjae	6.734	15.404	36.813	11.877	29.172
	Tan	5.605	17.953	33.751	15.686	27.005
	Jungnang	5.053	16.446	33.902	16.133	28.467
	Cheonggye	7.335	16.544	37.389	11.833	26.899
	Changreung	6.989	17.318	36.207	12.510	26.976
	Anyang	5.546	18.401	33.874	15.811	26.367
	Dorim	7.431	17.594	37.420	11.789	25.767

Table S7. Percent fluorescence response ($P_{(i,n)}$) values calculated by FRI of the five regions in the FEEM contour plot (for the FEEM-LowRes model) for wastewater effluents.

Region	Description	$P_{(1,n)}$	$P_{(2,n)}$	$P_{(3,n)}$	$P_{(4,n)}$	$P_{(5,n)}$
Seoul / Gyeonggi	S city, E treatment plant	6.230	19.496	37.373	12.756	24.145
	S city, T treatment plant	3.179	11.402	32.423	16.491	36.506
	S city, S treatment plant	4.916	14.971	36.536	12.993	30.585
	B city, G treatment plant	2.930	13.985	34.089	15.683	33.313
	B city, Y treatment plant	2.951	14.305	32.952	16.485	33.307
	P city, P treatment plant	9.794	14.168	32.760	13.658	29.621
	D city, D treatment plant	3.551	15.524	35.423	13.760	31.742
	U city, B treatment plant	3.012	14.880	35.699	14.316	32.093
South Jeolla	Y city	3.367	15.233	33.554	16.806	31.040
Gangwon	Y county	4.017	18.724	34.174	16.256	26.829
	Gangwon 1	7.754	22.274	35.152	12.172	22.648
	Gangwon 2	1.838	11.370	32.550	15.705	38.538
	Gangwon 3	4.117	18.881	35.318	14.330	27.354
	Gangwon 4	5.110	16.966	35.677	14.974	27.273
	Gangwon 5	1.637	9.618	33.454	14.686	40.605
	Gangwon 6	3.015	14.665	37.001	12.948	32.371
	Gangwon 7	5.950	20.737	34.737	14.559	24.017
	Gangwon 8	10.669	22.294	32.667	13.631	20.739
	Gangwon 9	3.841	14.460	32.966	17.045	31.688
	Gangwon 10	5.210	18.097	34.891	15.171	26.631
	Gangwon 11	4.824	17.142	34.525	15.727	27.782
	Gangwon 12	5.734	20.129	34.455	14.904	24.778
	Gangwon 13	2.523	12.286	34.683	14.047	36.461
	Gangwon 14	4.790	17.406	34.535	15.644	27.625
	Gangwon 15	3.988	18.278	30.890	19.604	27.240
	Gangwon 16	9.961	22.493	31.847	14.478	21.222
	Gangwon 17	4.839	19.140	33.242	16.402	26.378
	Gangwon 18	5.381	19.882	34.158	15.319	25.261
	Gangwon 19	2.522	9.736	47.172	9.209	31.361
	Gangwon 20	3.214	14.167	33.919	16.409	32.291

Table S8. RSS and AIC for the prediction of O₃ and •OH exposures.

Type	K*	O ₃ exposure		•OH exposure	
		RSS	AIC	RSS	AIC
FEEM-Free	3	5.4×10^{-5}	-141.6	5.4×10^{-20}	-556.2
FEEM-LowRes	8 (5)	4.6×10^{-5}	-132.0	4.5×10^{-20}	-548.4
FEEM-HighRes	86 (83)	4.2×10^{-5}	21.1	3.6×10^{-20}	-394.9
FEEM-FullRes	9,727 (9,724)	3.1×10^{-5}	19,299.6	3.7×10^{-20}	18,887.2

* K = the number of input variables (the values in parentheses are the number of FEEM variables)

Table S9. Second-order rate constants for the reactions of O₃ and •OH with MPs used in this study.

MP	k _{O3} (M ⁻¹ s ⁻¹)	k _{•OH} (M ⁻¹ s ⁻¹)	Reference
Atrazine	4.0	2.7×10^9	S1
Ibuprofen	9.6	7.4×10^9	S2
Caffeine	6.5×10^2	5.9×10^9	S3
Carbamazepine	3.0×10^5	8.8×10^9	S2

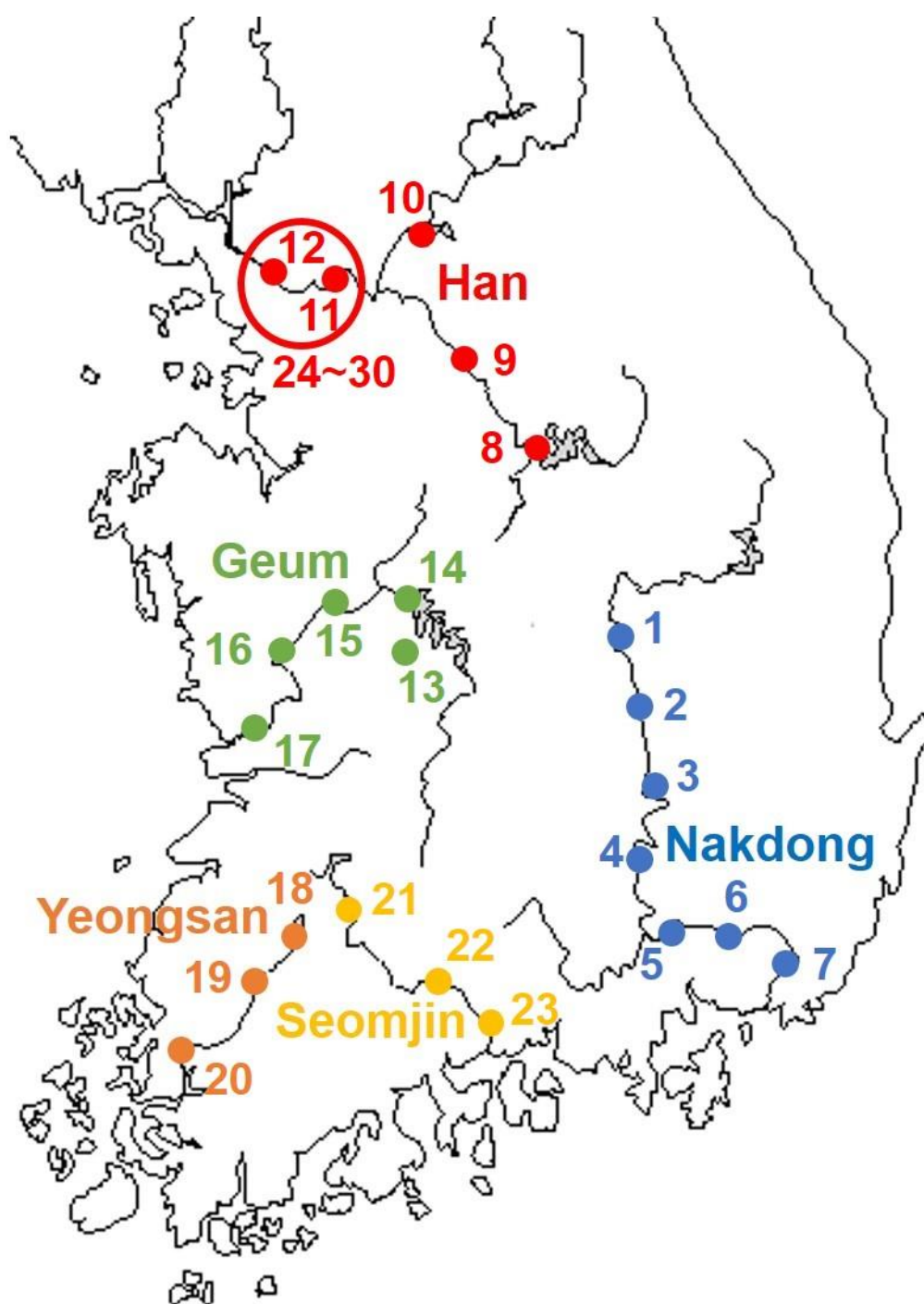


Figure S1. Sampling locations of natural waters across the four major rivers in South Korea. Note that samples for tributaries of the Han River were taken in the red circled region. Refer to Table S2 for the numbers of natural water sampling locations.

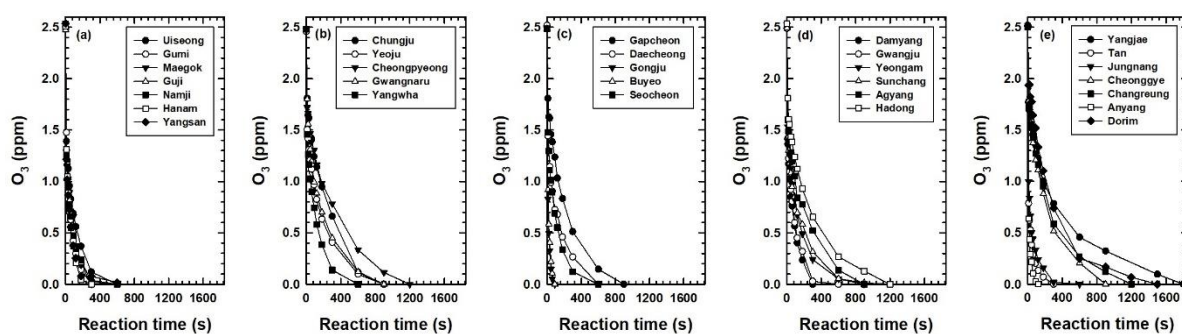


Figure S2. Time–concentration profiles of O_3 during the ozonation of natural waters, (a) Nakdong [1–7], (b) Han [8–9], (c) Geum [13–17], (d) Yeongsan and Sumjin Rivers [18–23], and (e) tributaries of the Han River [24–30]. Numbering in brackets corresponds to sampling locations as designated in Table S2.

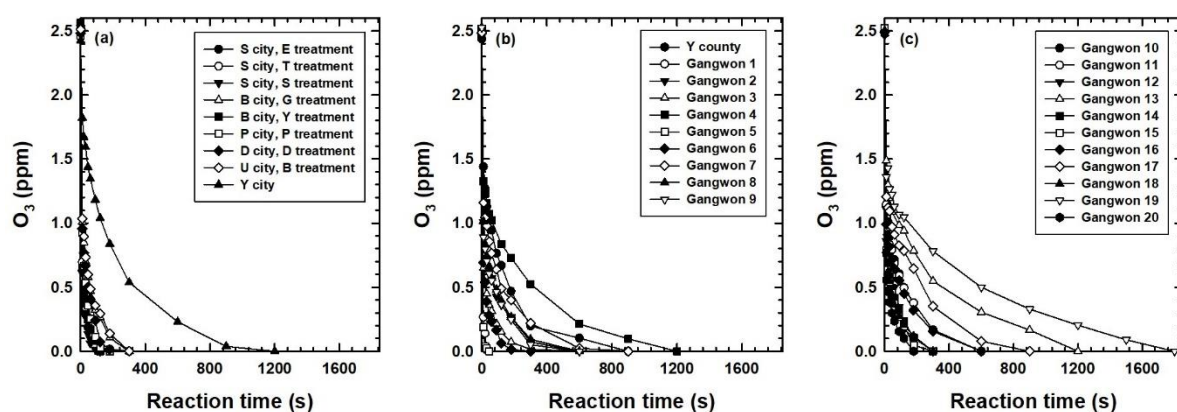


Figure S3. Time–concentration profiles of O_3 during the ozonation of wastewater effluents from WWTPs in (a) Seoul, and Gyeonggi and South Jeolla provinces [1–9], and (b, c) Gangwon province [10–30]. Numbering in brackets corresponds to sampling locations as designated in Table S3.

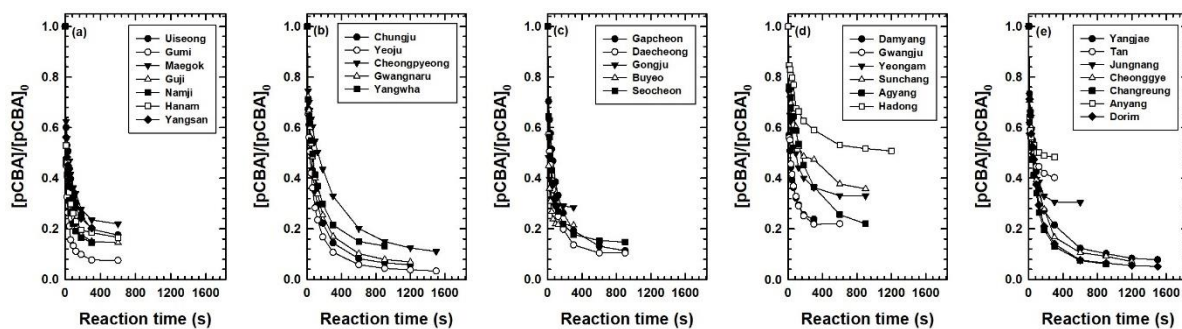


Figure S4. Time–concentration profiles of *p*CBA during the ozonation of natural waters, (a) Nakdong, (b) Han, (c) Geum, (d) Yeongsan and Sumjin Rivers, and (e) tributaries of the Han River.

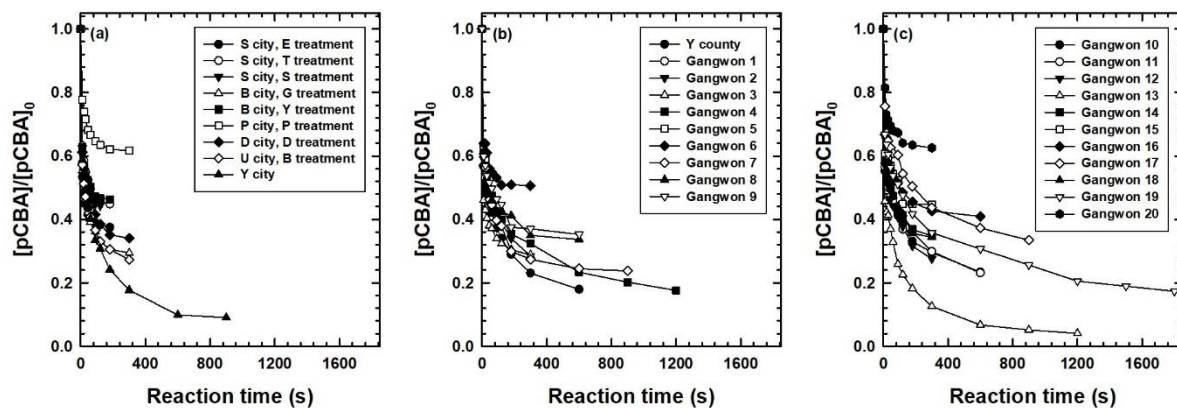


Figure S5. Time–concentration profiles of *p*CBA during the ozonation of wastewater effluents from WWTPs in (a) Seoul, and Gyeonggi and South Jeolla provinces, and (b, c) Gangwon province.

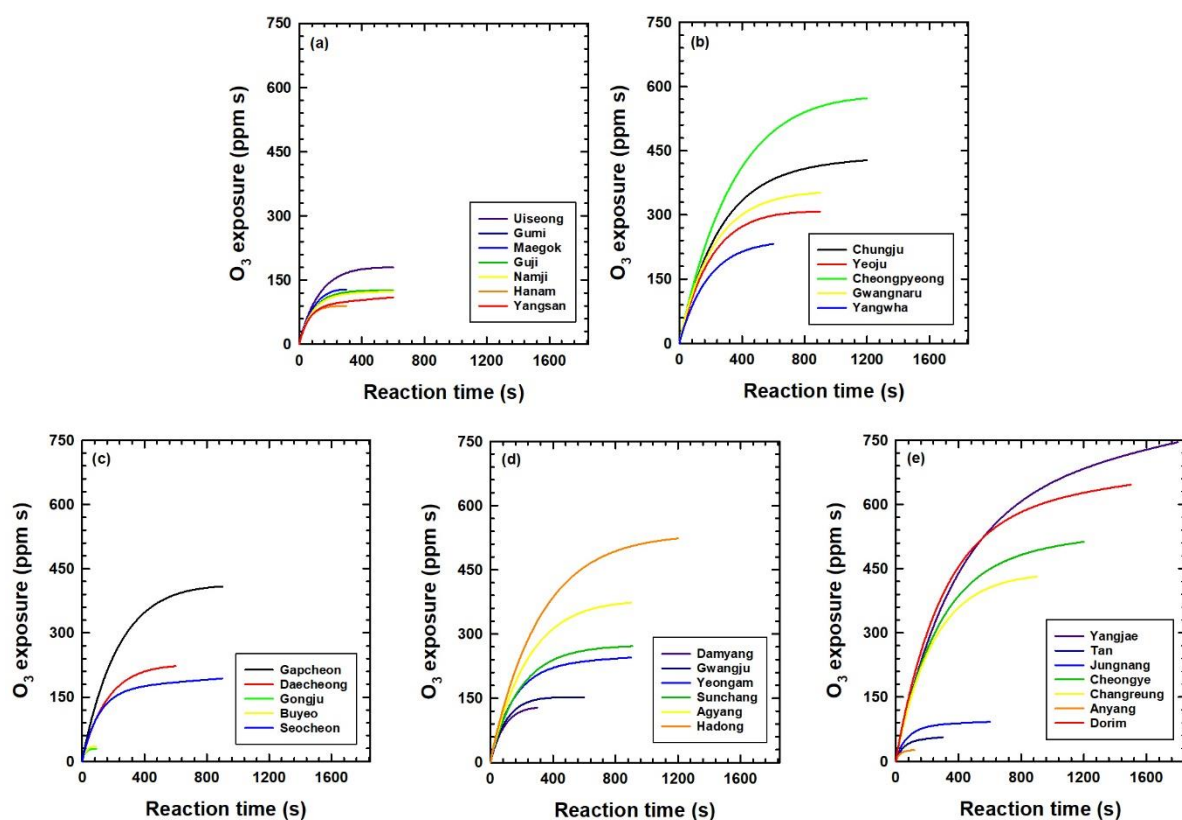


Figure S6. Time-dependent O_3 exposure during the ozonation of natural waters, (a) Nakdong, (b) Han, (c) Geum, (d) Yeongsan and Sumjin Rivers, and (e) tributaries of the Han River.

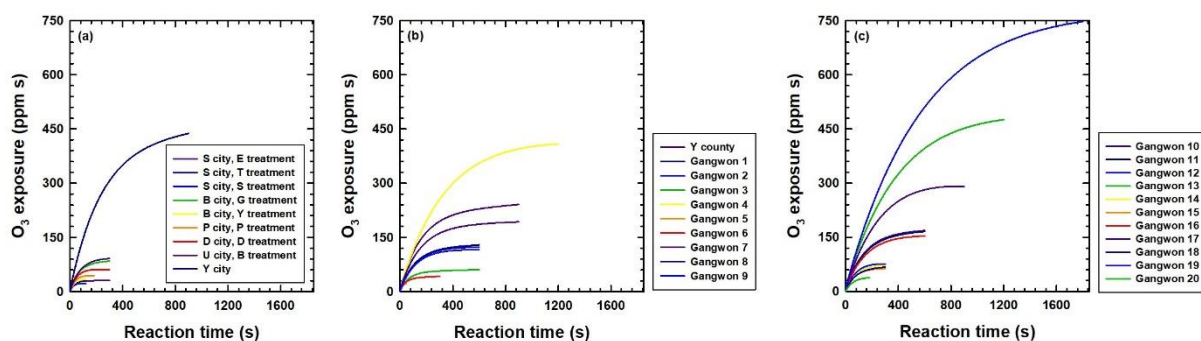


Figure S7. Time-dependent O_3 exposure during the ozonation of wastewater effluents from WWTPs in (a) Seoul, and Gyeonggi and South Jeolla provinces, and (b, c) Gangwon province.

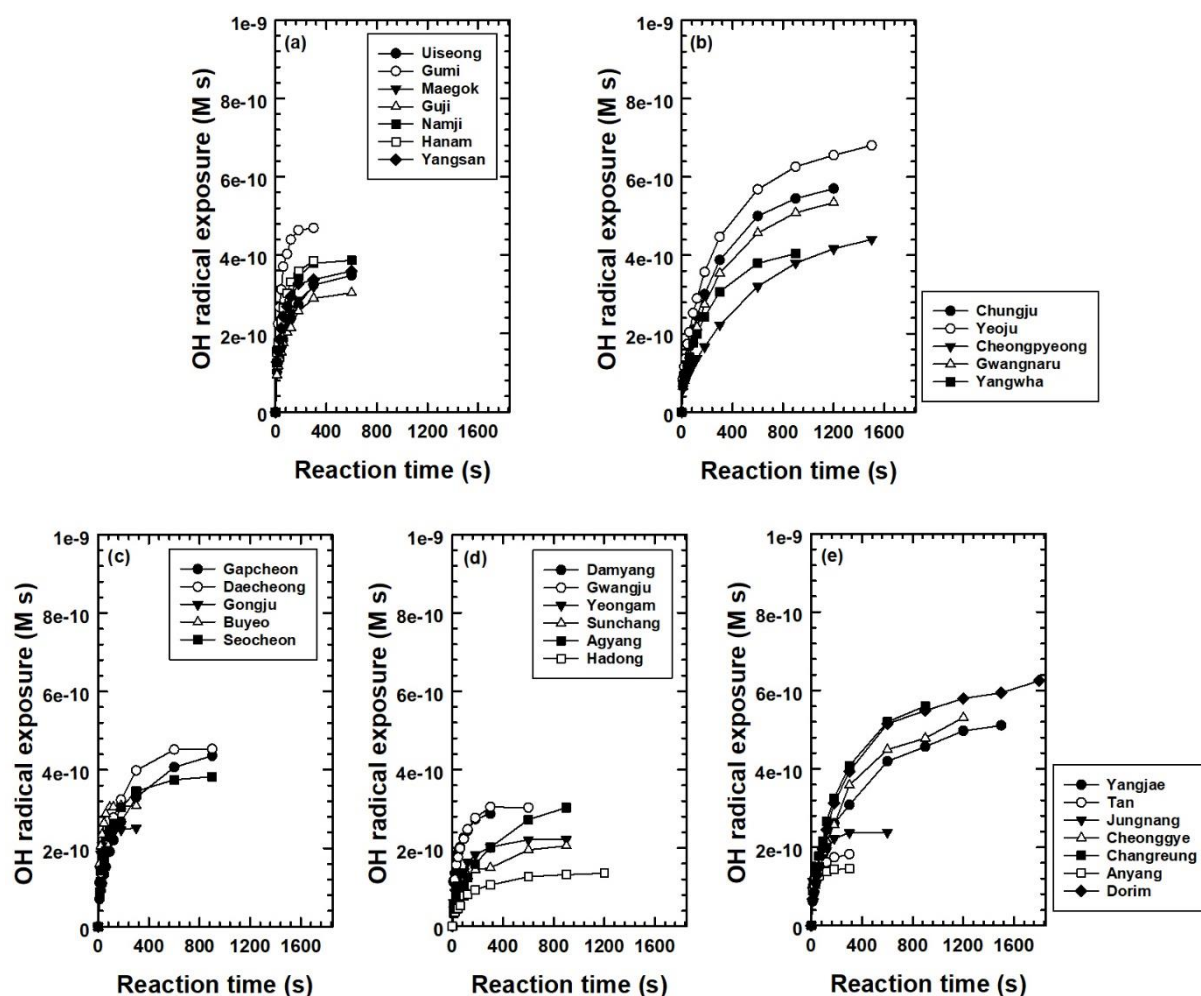


Figure S8. Time-dependent $\bullet\text{OH}$ exposure during the ozonation of natural waters, (a) Nakdong, (b) Han, (c) Geum, (d) Yeongsan and Sumjin Rivers, and (e) tributaries of the Han River.

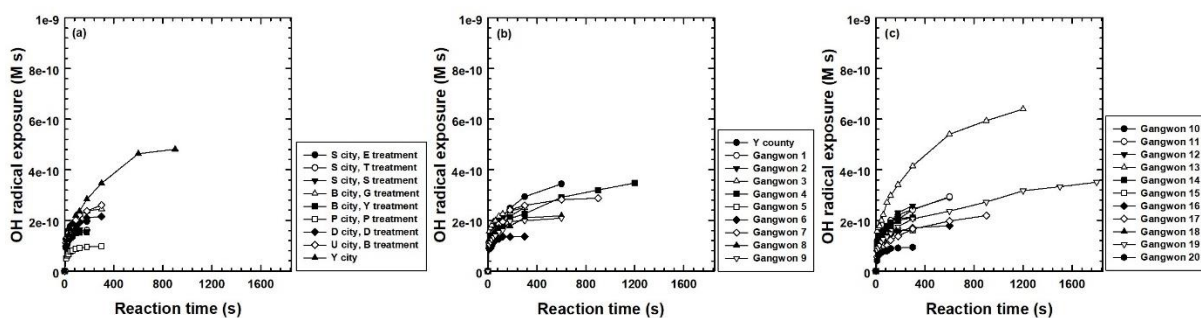


Figure S9. Time-dependent $\bullet\text{OH}$ exposure during the ozonation of wastewater effluents from WWTPs in (a) Seoul, and Gyeonggi and South Jeolla provinces, and (b, c) Gangwon province.

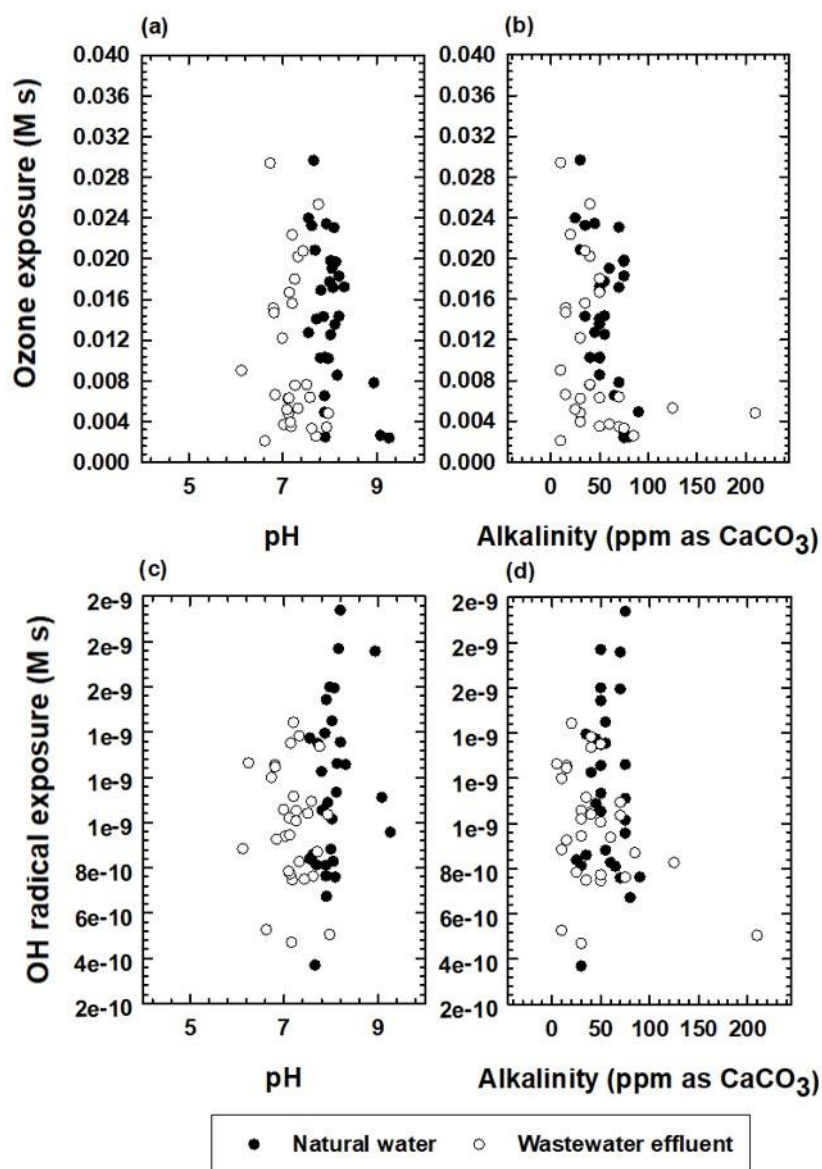


Figure S10. O₃ exposures as functions of (a) pH, and (b) alkalinity. •OH exposures as functions of (c) pH, and (d) alkalinity ([O₃]₀ = 2.5 mg/L).

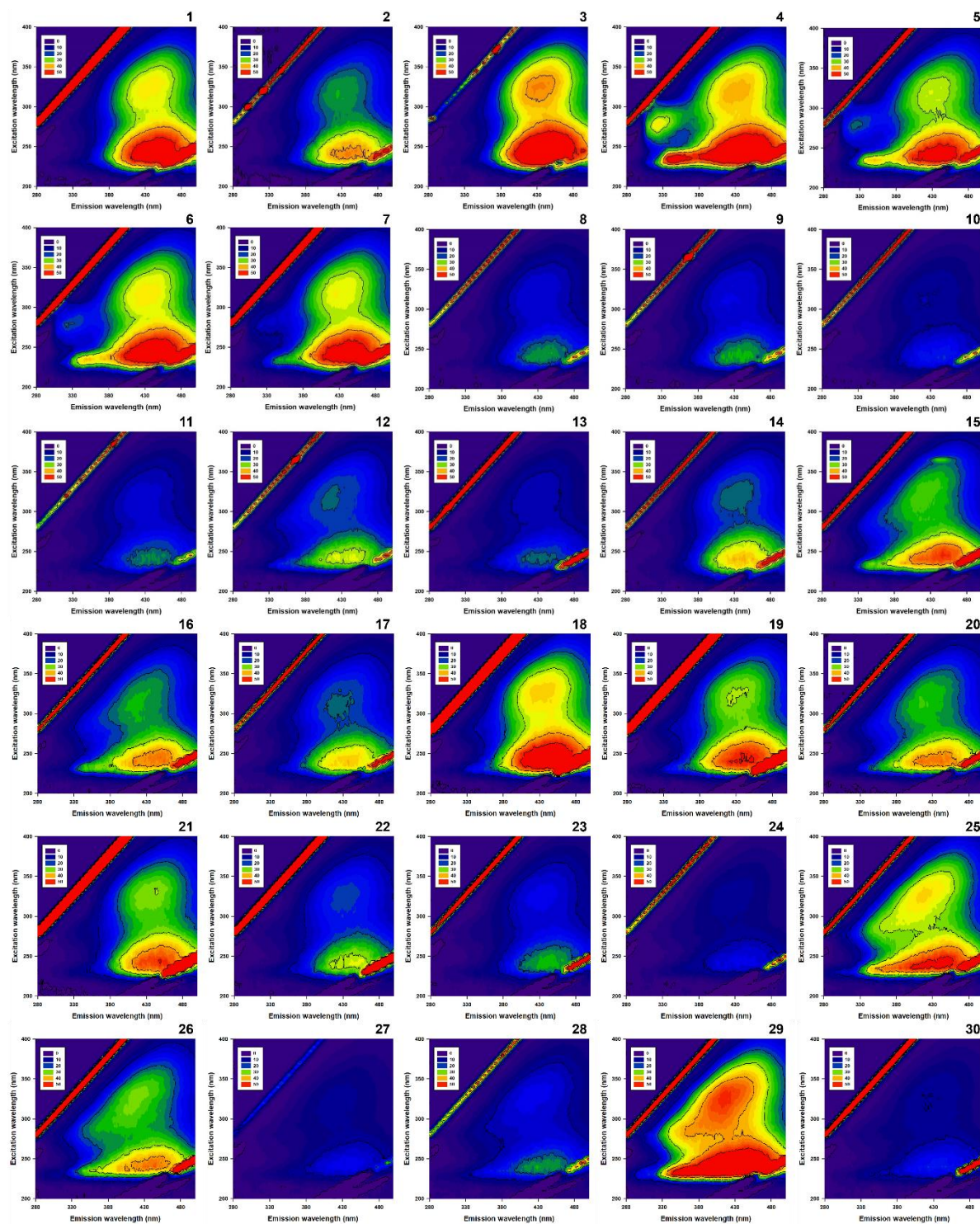


Figure S11. FEEM contour plots of natural waters. Numbering corresponds to sampling locations as designated in Table S2.

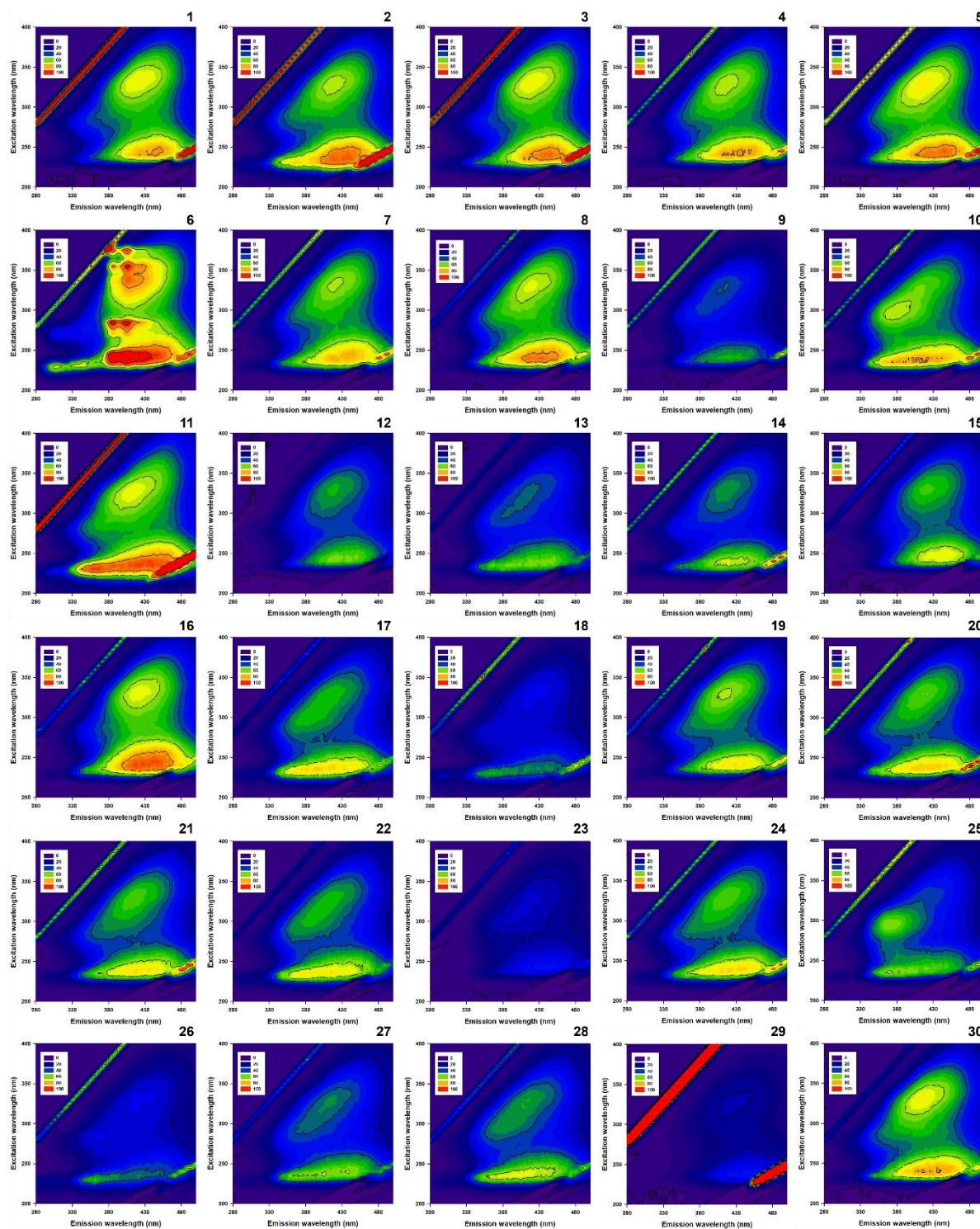


Figure S12. FEEM contour plots of wastewater effluents. Numbering corresponds to sampling locations as designated in Table S3.

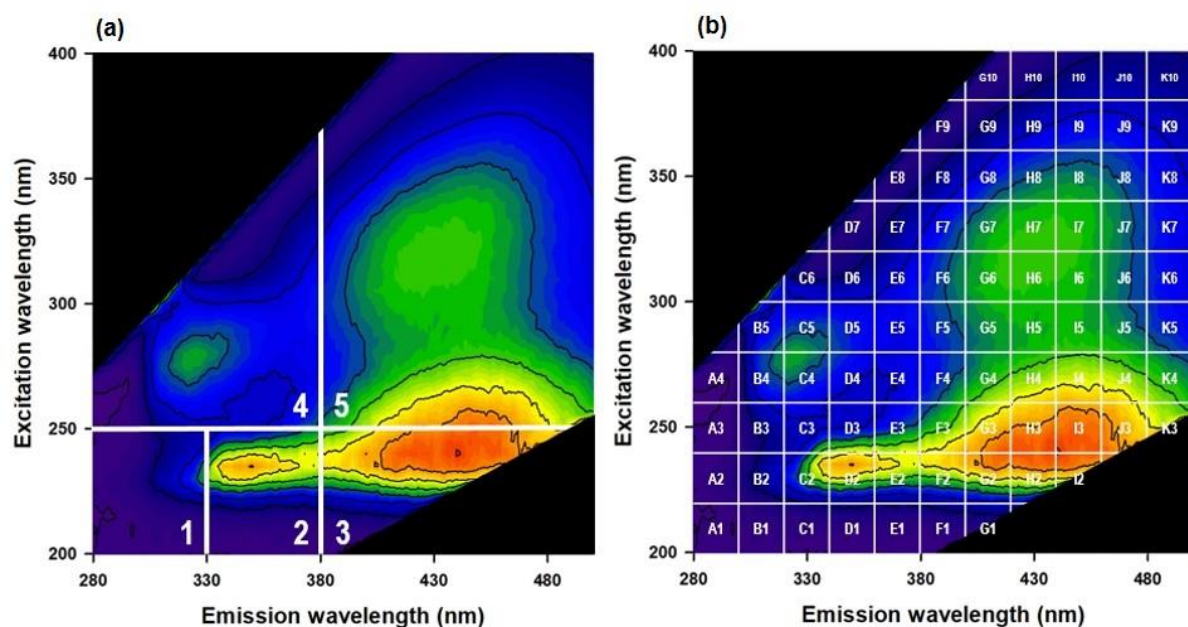


Figure S13. FRI division for (a) FEEM-LowRes, and (b) FEEM-HighRes models (the FEEM contour plot of Figure 3a was used as an example).

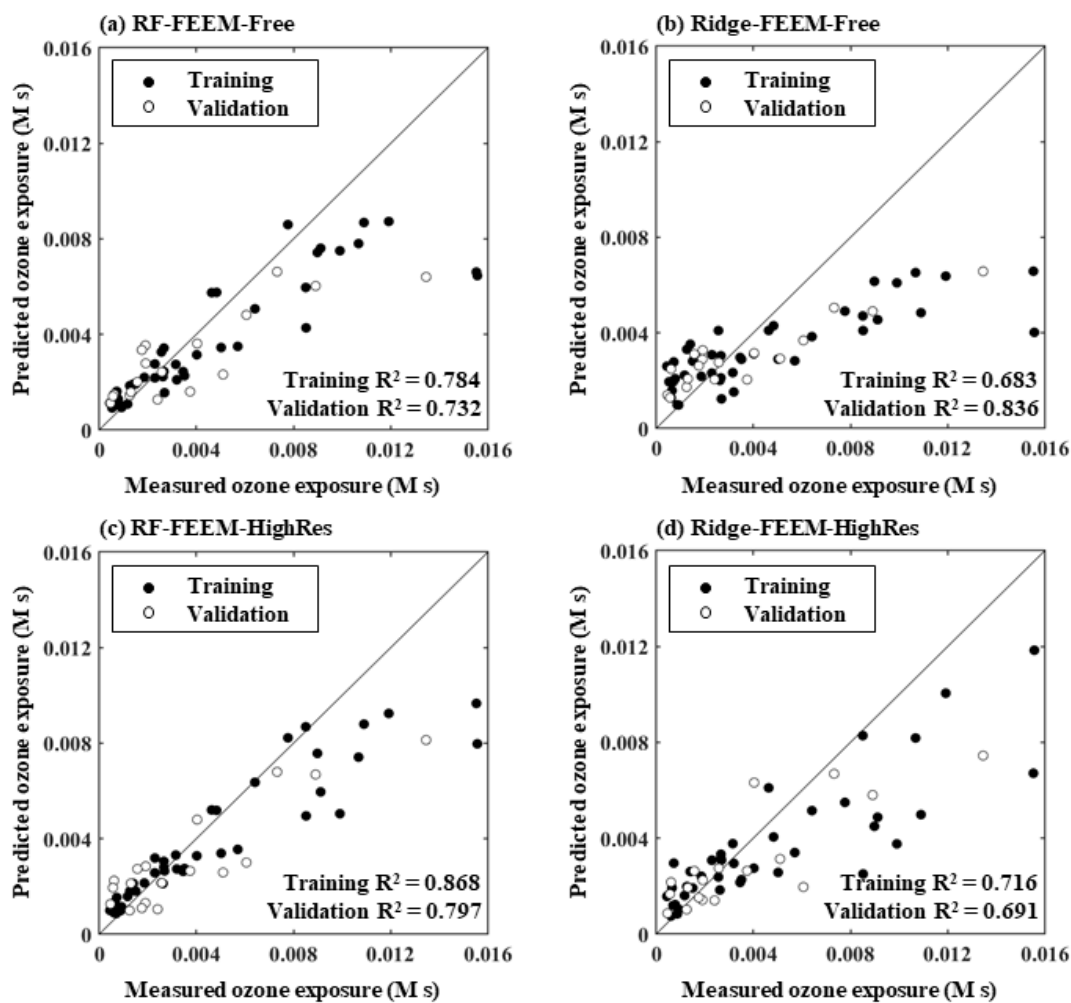


Figure S14. Plots of measured versus predicted O₃ exposures for RF and ridge regression models, (a) RF-FEEM-Free, (b) Ridge-FEEM-Free, (c) RF-FEEM-HighRes, and (d) Ridge-FEEM-HighRes ([O₃]₀ = 2.5 mg/L).

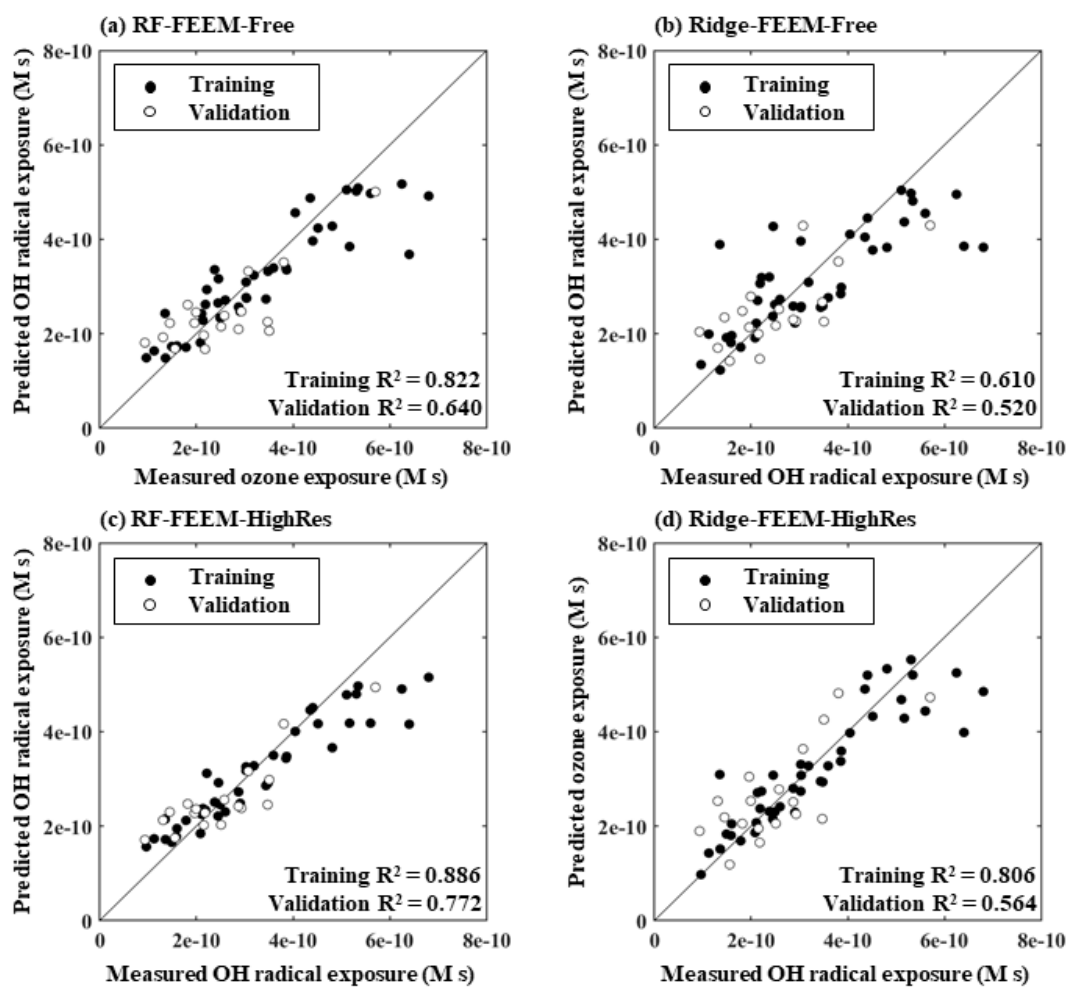


Figure S15. Plots of measured versus predicted $\bullet\text{OH}$ exposures for RF and ridge regression model, (a) RF-FEEM-Free, (b) Ridge-FEEM-Free, (c) RF-FEEM-HighRes, and (d) Ridge-FEEM-HighRes ($[\text{O}_3]_0 = 2.5 \text{ mg/L}$).

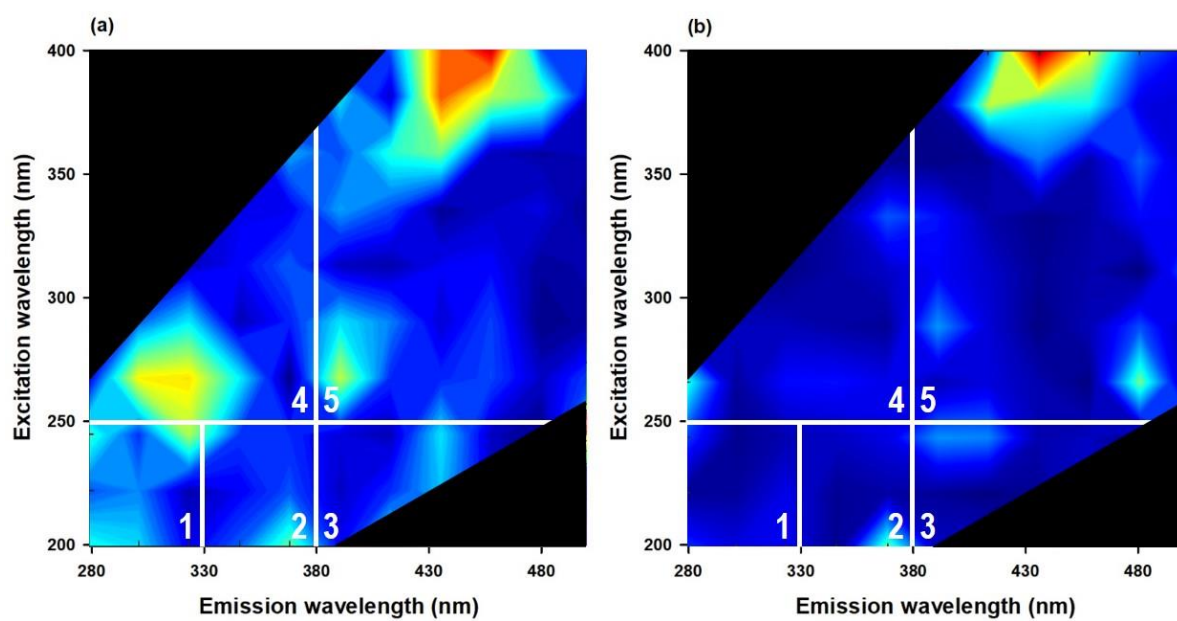


Figure S16. Variable importance (VI) for the prediction of (a) O_3 exposure, and (b) $\bullet OH$ exposure by the FEEM-HighRes model.

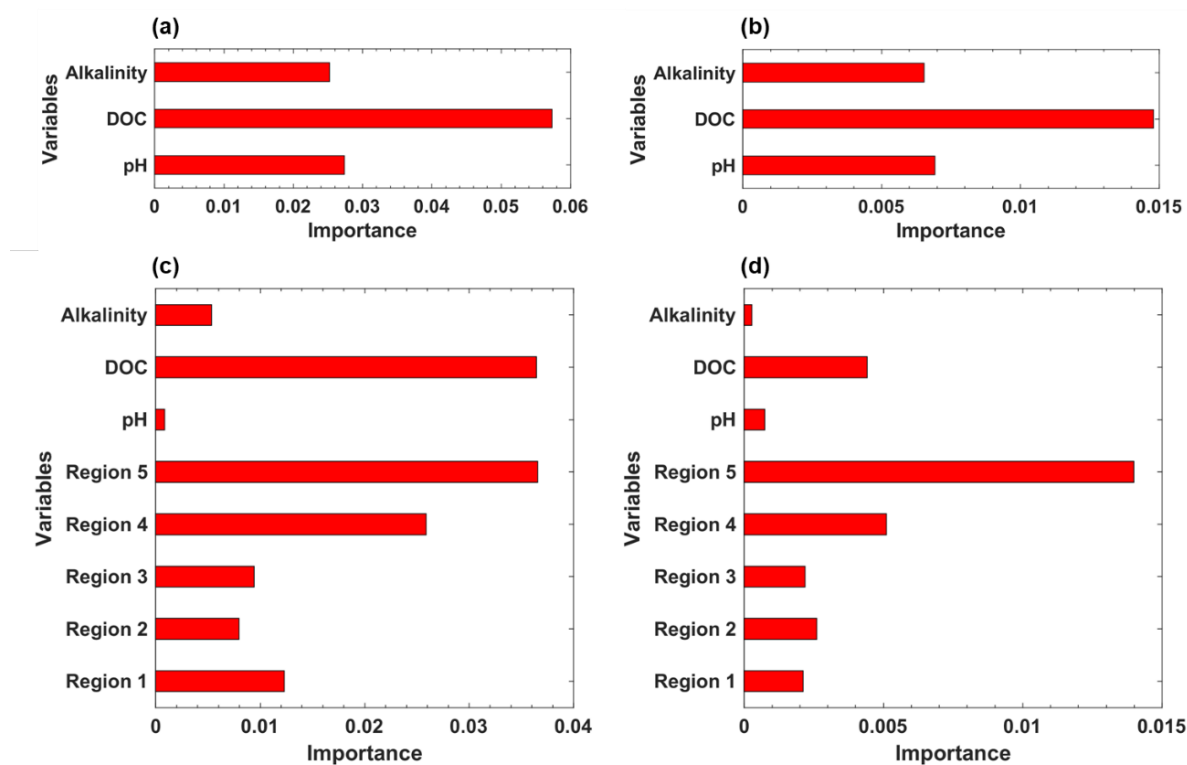
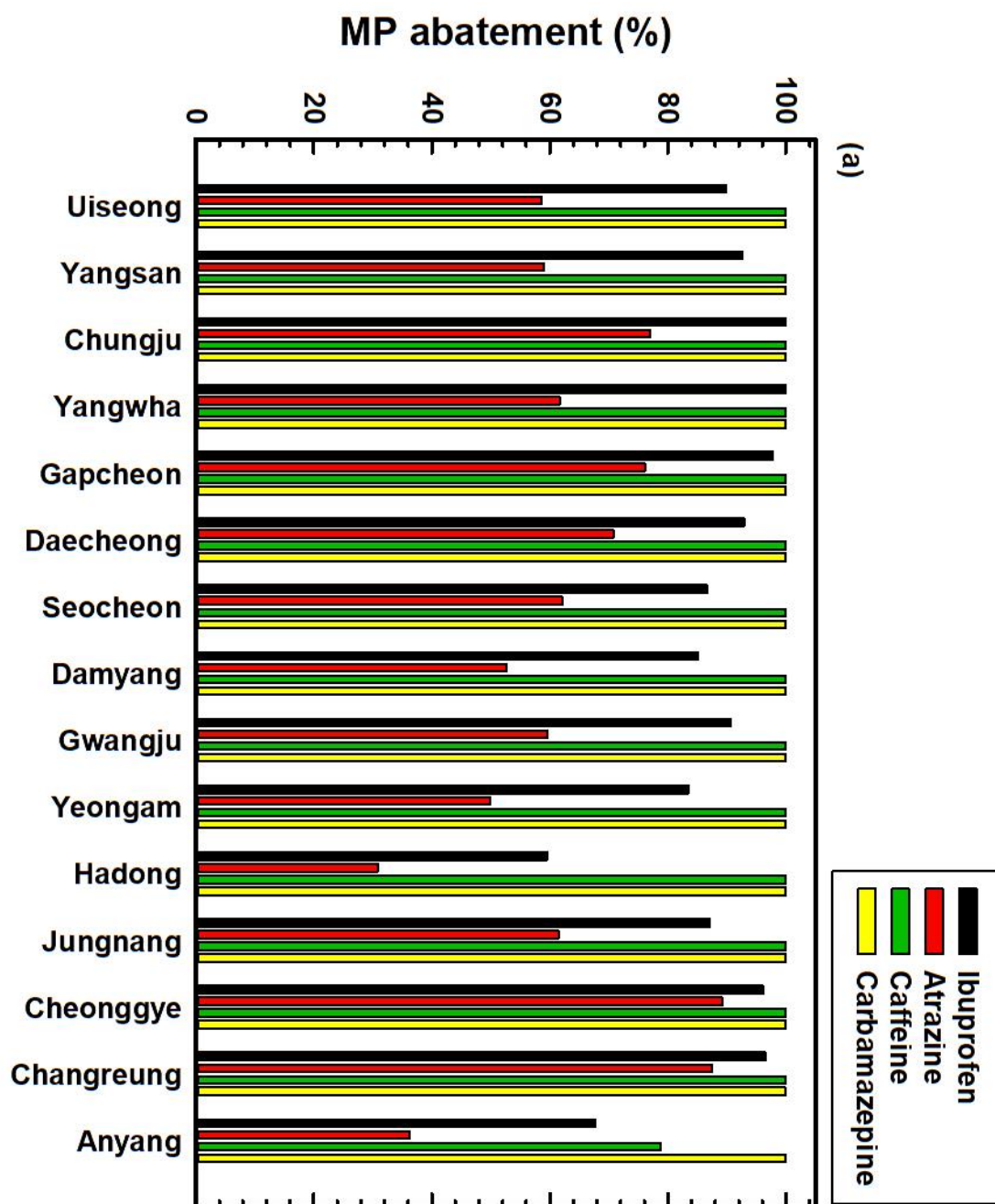


Figure S17. Variable importance (VI) for the prediction of (a) O₃ exposure and (b) •OH exposure by the FEEM-Free model, (c) O₃ exposure and (d) •OH exposure by the FEEM-HighRes model.



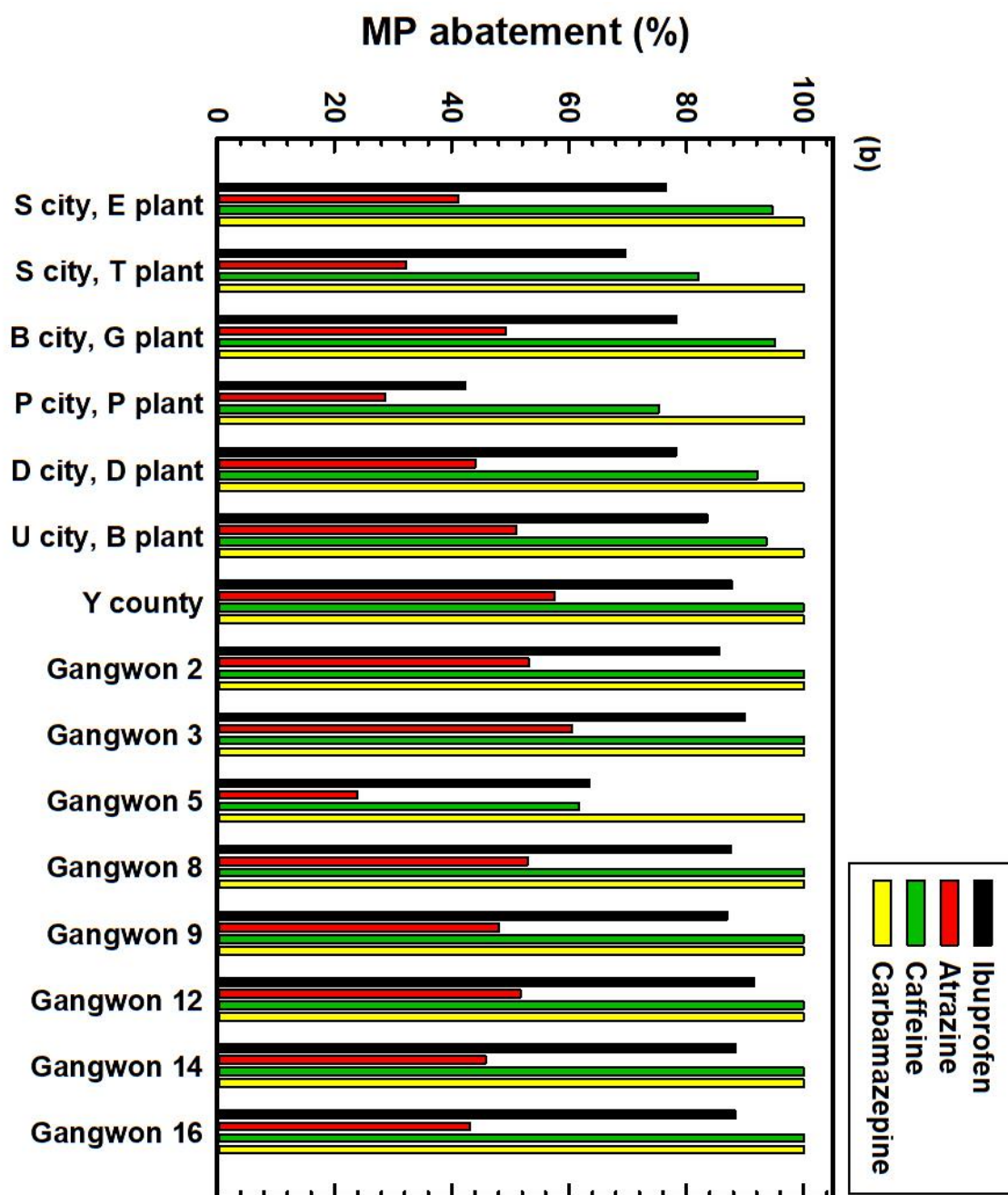


Figure S18. Percent removal of MPs in (a) natural water samples, and (b) wastewater effluent samples (15 natural water and wastewater effluent samples, each, i.e., total 30 samples).

Supporting References

- S1. Acero, J. L.; Stemmler, K.; von Gunten, U. Degradation kinetics of atrazine and its degradation products with ozone and OH radicals: A predictive tool for drinking water treatment. *Environ. Sci. Technol.* **2000**, *34* (4), 591-597. DOI: 10.1021/es990724e.
- S2. Huber, M. M.; Canonica, S.; Park, G.-Y.; von Gunten, U. Oxidation of pharmaceuticals during ozonation and advanced oxidation processes. *Environ. Sci. Technol.* **2003**, *37* (5), 1016-1024. DOI: 10.1021/es025896h.
- S3. Broséus, R.; Vincent, S.; Aboulfadl, K.; Daneshvar, A.; Sauvé, S.; Barbeau, B.; Prévost, M. Ozone oxidation of pharmaceuticals, endocrine disruptors and pesticides during drinking water treatment. *Water Res.* **2009**, *43* (18), 4707-4717. DOI: 10.1016/j.watres.2009.07.031.