

Supplement of Observed changes in the temperature dependence response of surface ozone under NO_x reductions over Germany

Noelia Otero, Henning W. Rust and Tim Butler

Modeling ozone production rates with GAMs

Generalized Additive Models (GAMs) (Hastie and Tibshinari 1990; Wood 2006) are useful tools to examine complex non-linear relationships and have been previously applied to model air pollutants (Barnpadimos et al. 2011; Boleti, Hueglin, and Takahama 2019; Carslaw, Beevers, and Tate 2007; Jackson et al. 2009). We have used GAMs to model O₃ production rates (ΔO_3) as a function of key variables that influence O₃ production. GAMs are extensions of the generalized linear model (McCullagh and Nelder 1989) that work under the assumption that there is an additive effect between the response variable and the explanatory variables (covariates). Generalized linear models allow for response distributions other than the Normal distribution, and for a degree of non-linearity in the model structure (McCullagh and Nelder 1989). The basic form of GLM is represented as:

$$g(\mu_i) = X_i\beta \quad (1)$$

where $\mu_i = E[Y_i]$, g is a monotonic function, X_i is the i^{th} row of X (model matrix) and β is a vector of unknown parameters. GLM assumes that Y_i are independent and $Y_i \sim$ some exponential family distribution (for more details see McCullagh and Nelder 1989; Wood 2006).

As stated in the manuscript (see section 3.3), GAM establishes a relationship between the response and a sum of smooth functions of the covariates through a link function (Hastie and Tibshinari 1990; Wood 2006). Thin plate regression splines were used as smoothers to describe a nonlinear relationship between the response and the covariates (Wood 2006). In addition, GAMs allow to model interactions created between covariates with different smoothers (or degrees of smoothness) assumed for each covariate (Wood 2006; Pedersen et al. 2019). Here, we introduced interactions terms using tensor products to represent the the interacting effects of two covariates (e.g. temperature-NO_x) on the response (ΔO_3). For a general overview of GAM we refer to Hastie and Tibshinari (1990) and Wood (2006).

All calculations were carried out using the statistical software R (R Development Core Team 2018) with the mgcv package (Wood 2011).

Model selection

A set of covariates were used to build the GAMs: temperature, NO_x, VPD, O₃ concentrations from the previous hour ($C_{O_3}(t-1)$), boundary layer height growth rate (ΔBLH) and the MDA8 concentrations from the previous day ($C_{MDA8}(t-24)$). A forward selection process was used to select the covariates that better explain the ΔO_3 . During the selection procedure, the interactions between two influencing covariates are also included in order to represent physical processes such as dry deposition, represented by the interaction between VPD and $C_{O_3}(t-1)$, and mixing processes captured by the interaction term between ΔBLH and $C_{MDA8}(t-24)$.

The selection process can be summarised as follows:

1. We first start with a baseline model that included the nonlinear relationship between NO_x and temperature as follows:

$$\Delta O_3 = f(T, NO_x) \quad (2)$$

where $f(T, NO_x)$ represents the interaction between temperature (T) and NO_x concentrations and it is included as a tensor product (Wood 2017). Observing the skewness of the NO_x data led us to introduce a modification in the baseline model using a log transformation of NO_x .

2. We successively add further covariates and/or interactions that can improve the model performance.
3. The deviance explained and the Akaike information criterion (AIC) (Akaike 1974) are calculated in each step.
4. The GAM with the lowest AIC is selected as the best model.

We applied this procedure separately for each station and period, namely GAM-P1 for the first period (1999-2008) and GAM-P2 for the second period (2009-2019). Our goal with this process is to define a common model well defined across all of the stations (i.e. same structure in terms of covariates). Figure S1 shows the models built at each step (i.e. adding the covariates and interactions) during the selection process for the urban station in Berlin during the first period 1999-2008. It can be observed that the model performance considerably improves when adding the covariates and the complexity (i.e. more interaction terms).

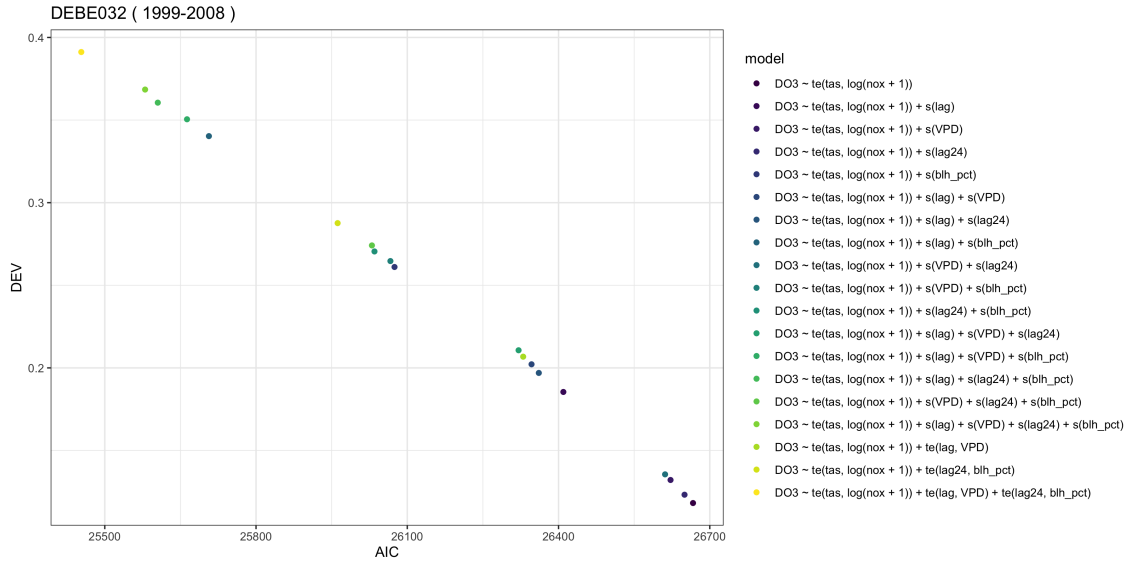


Figure S1. AIC and deviance explained (DEV) for each model used during the stepwise process at one rural station during the first period 1999-2008.

We obtained similar results for most of stations, which led us to select the best model with the following structure that includes three interaction terms:

$$\Delta O_3 = T * NO_x + VPD * C_{O_3}(t - 1) + \Delta BLH * C_{MDA8}(t - 24) \quad (3)$$

The model performance was assessed through standard diagnostic plots: QQ plots of the deviance residuals, scatter plots of the residuals against the fitted values, histogram of residuals and scatter plots the response against the fitted values (Wood 2006). In general the diagnostic plot did not show concerning patterns in the residuals. As an example, Fig. S2 shows the standard plots to check the model assumptions obtained by the function *gam.check()* (Wood 2011).

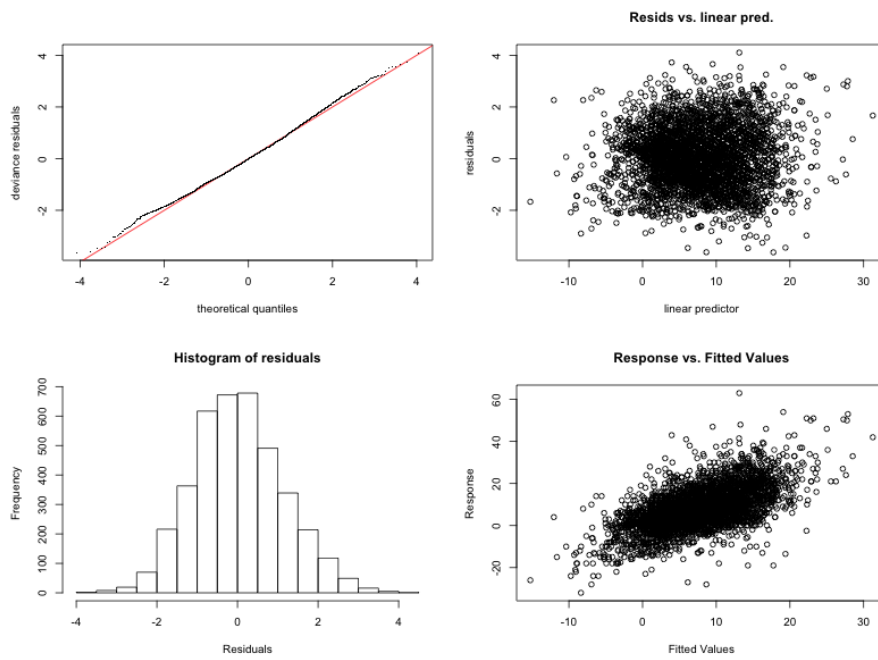


Figure S2. Diagnostic plots for the Berlin urban station (DEBE034) for the period 1999-2008: QQ-plot of residuals, linear predictor vs. residuals, the histogram of residuals and the plot of fitted values vs. response.

List of figures

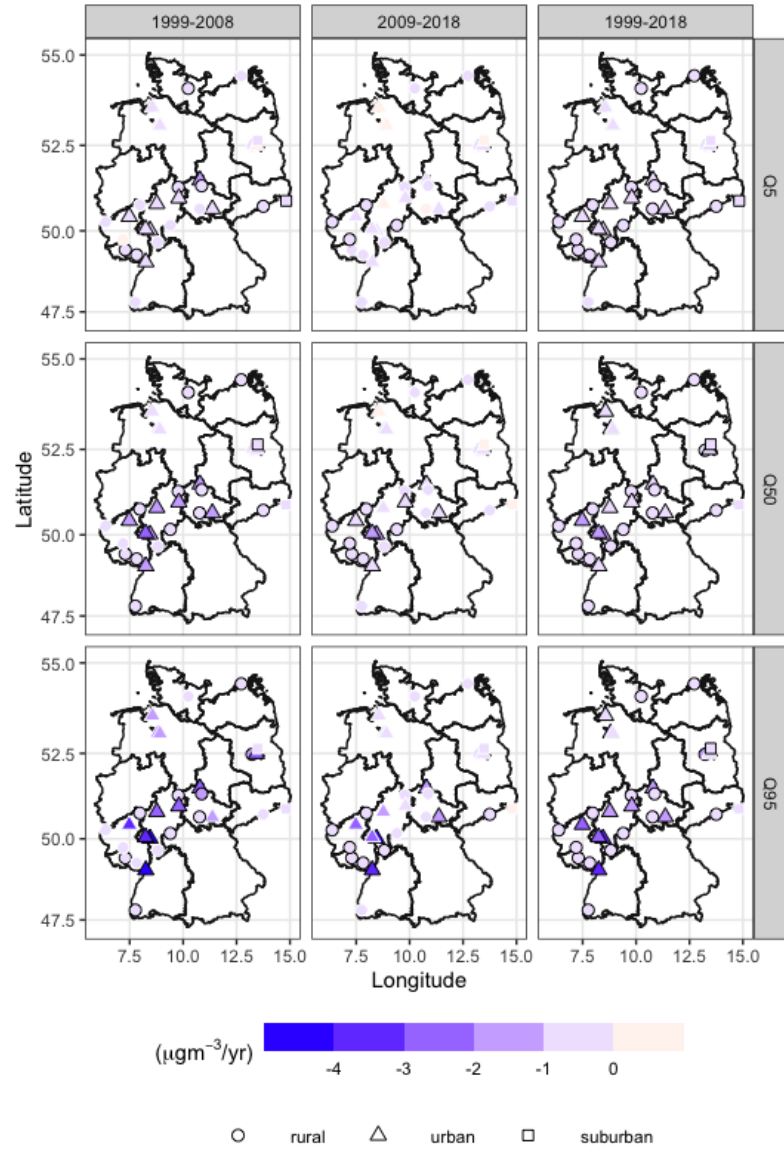


Figure S3. Spatial distribution of trends calculated separately for each station and period, 1999-2008, 2009-2018, and the complete period 1999-2018. Bold black circles represent significant trends.

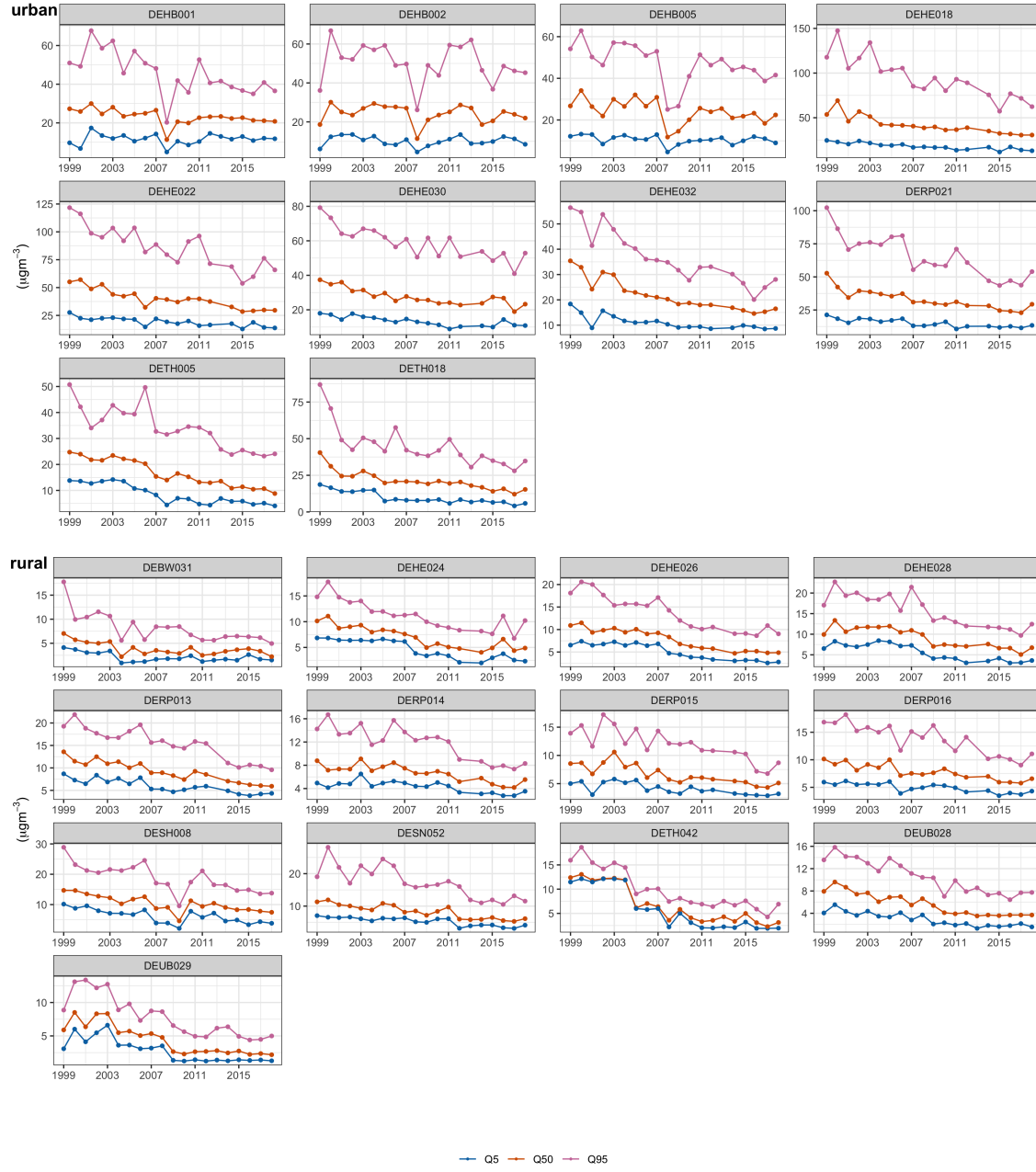


Figure S4. Time series of the annual 5th, 50th, 95th, percentiles for the rest of the urban and rural stations not presented in the main text

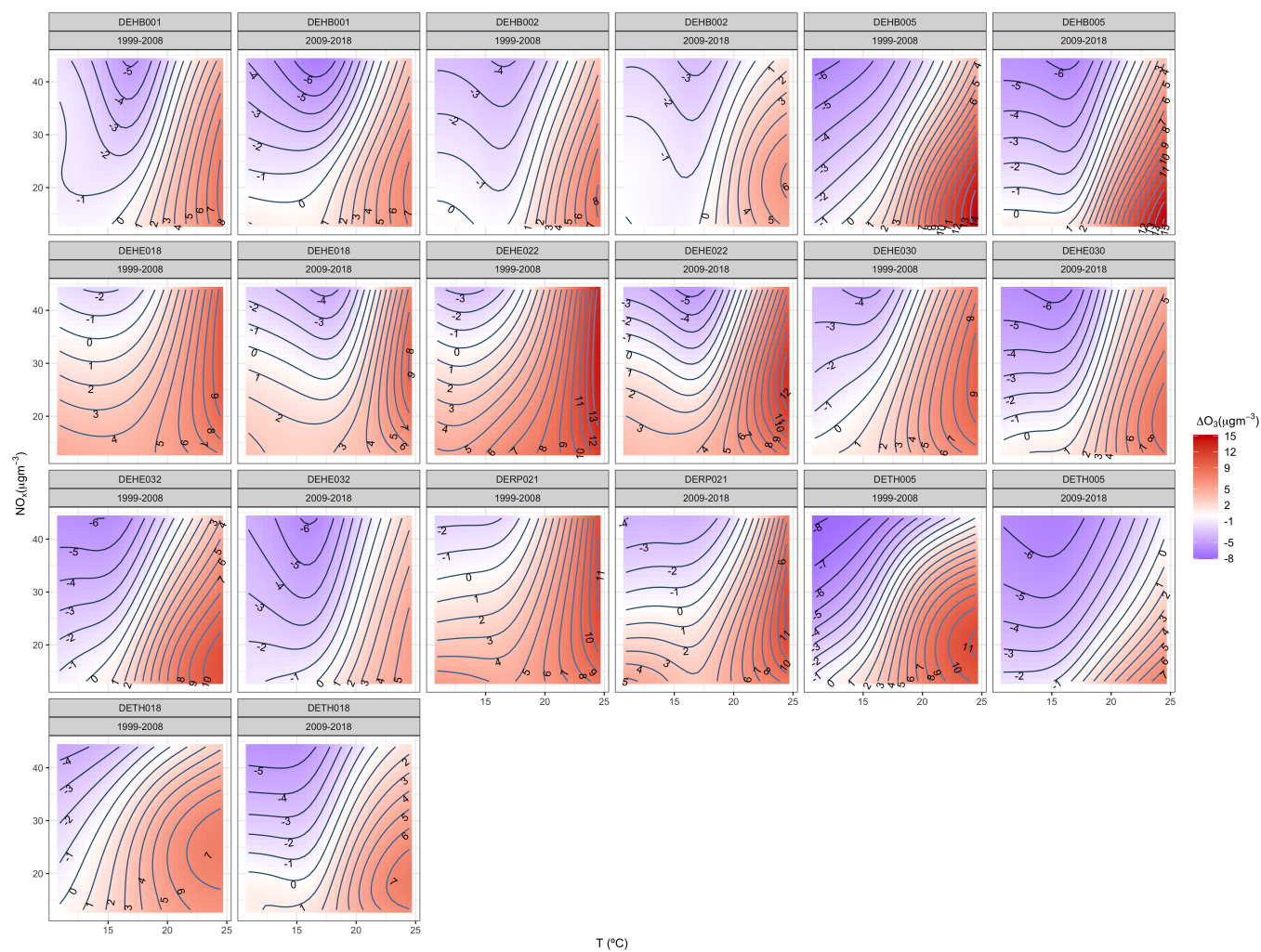


Figure S5. Countour plots obtained for the interaction term temperature and NO_x from each GAMs built separately at each urban station and corresponding period.

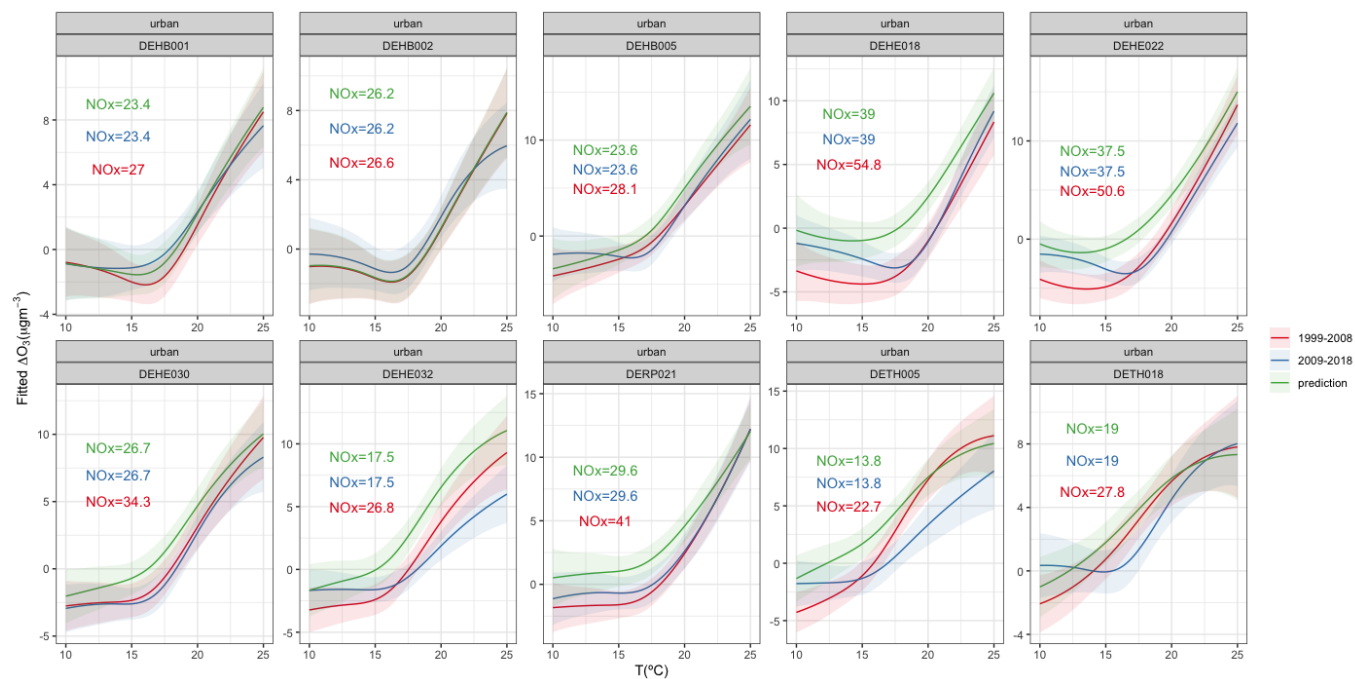


Figure S6. Estimated regression lines for urban stations of the ozone response to temperature while holding NOx concentrations constant (mean values for each period and station). Red line correspond to the prediction for the first period 1999-2008 (GAM-P1), blue line corresponds to the second period 2009-2018 and green line corresponds to the projected response using GAM-P1 with mean NOx conditions of the second period. Shaded bands represent the pointwise 95% confidence interval.

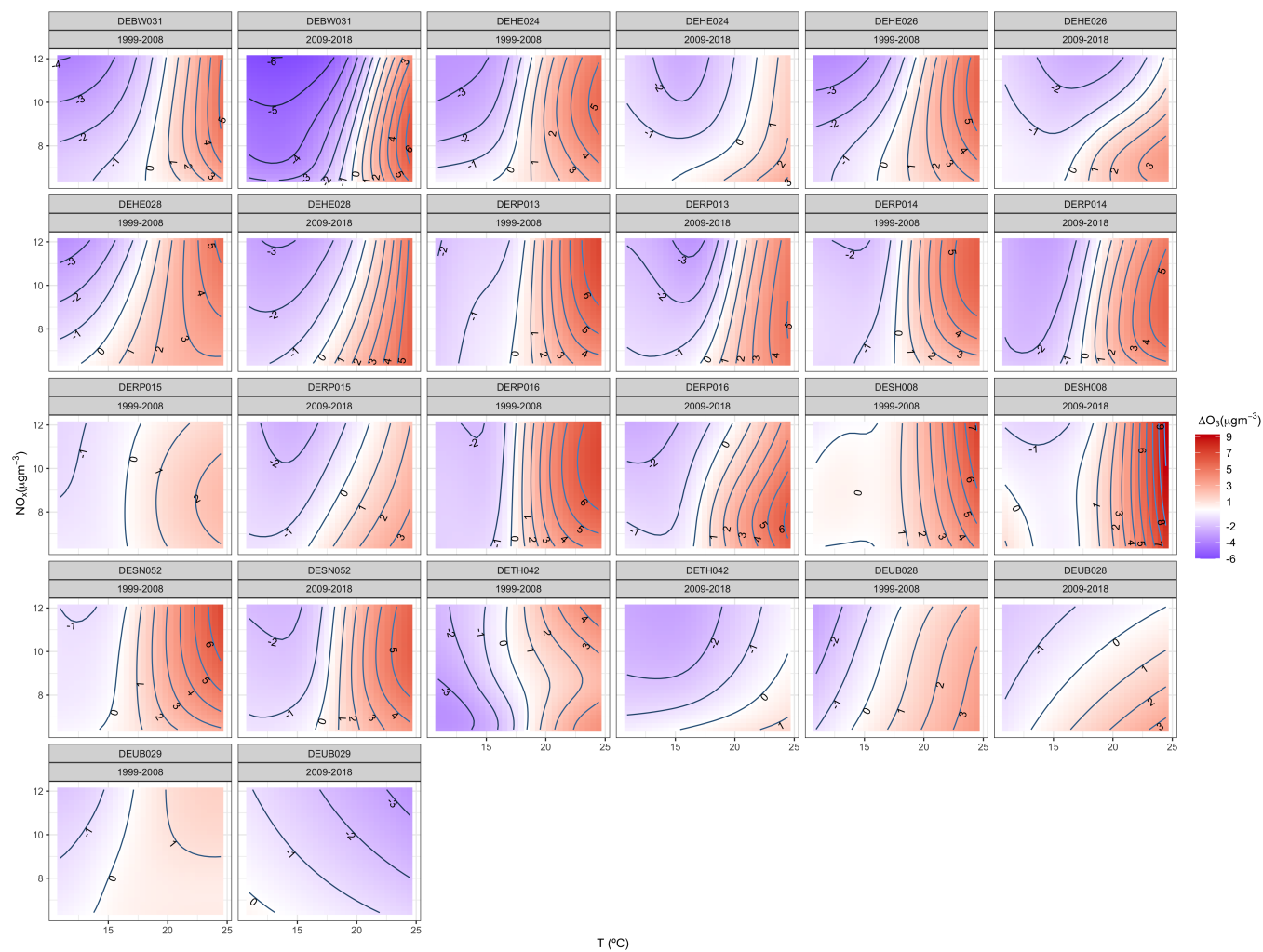


Figure S7. As Fig. S5, but for rural stations.

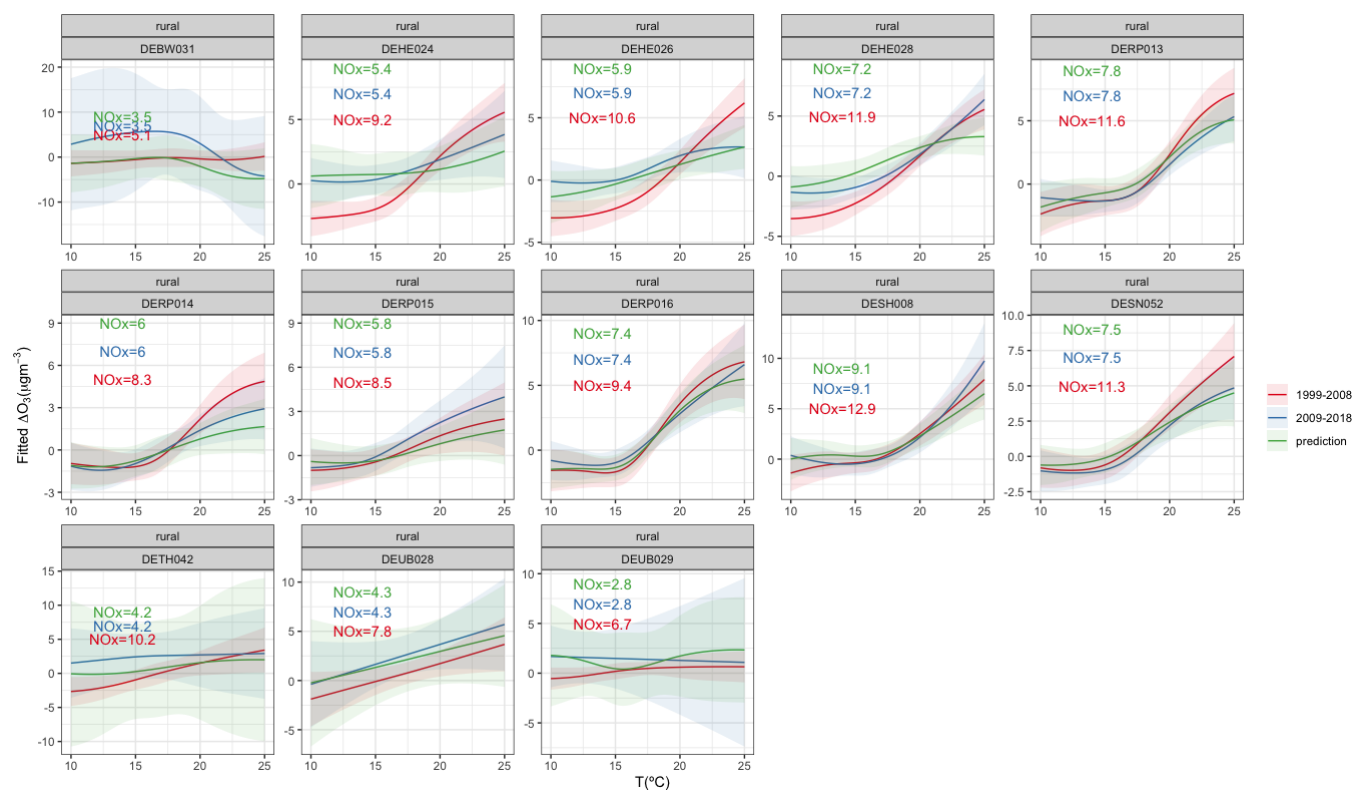


Figure S8. As Fig. S6, but for rural stations.

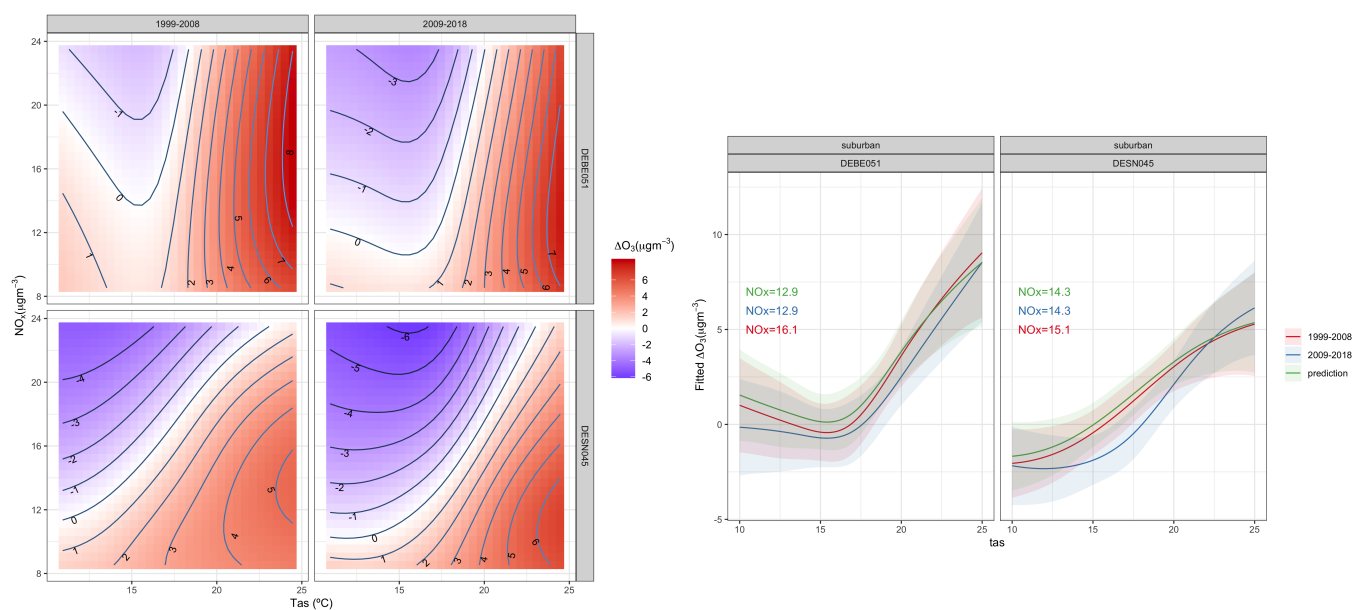


Figure S9. As Fig. 5 in the main text., but for the suburban stations.

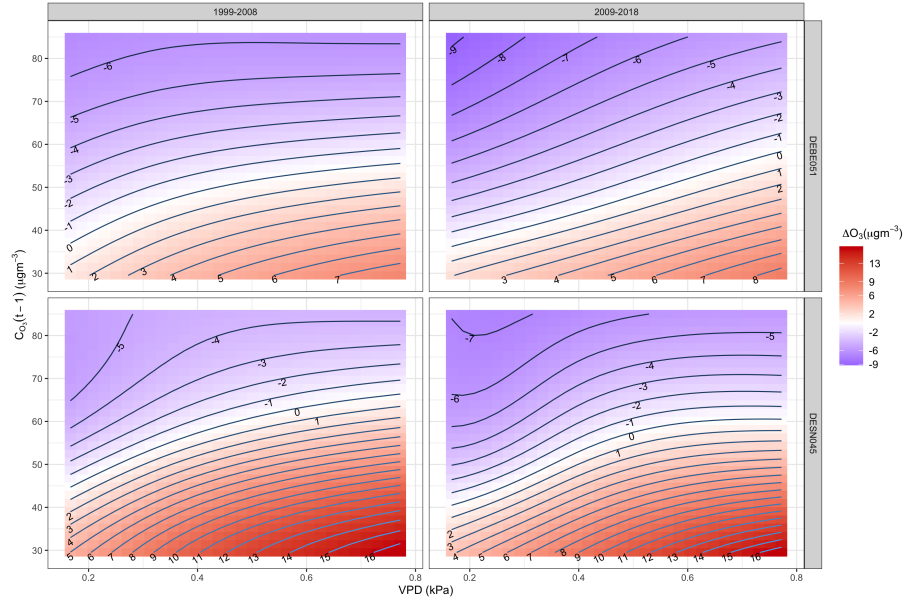


Figure S10. Contour plot for the interaction $VPD-CO_3(t-1)$ at the suburban stations for the first period 1999-2008 and second period 2009-2018.

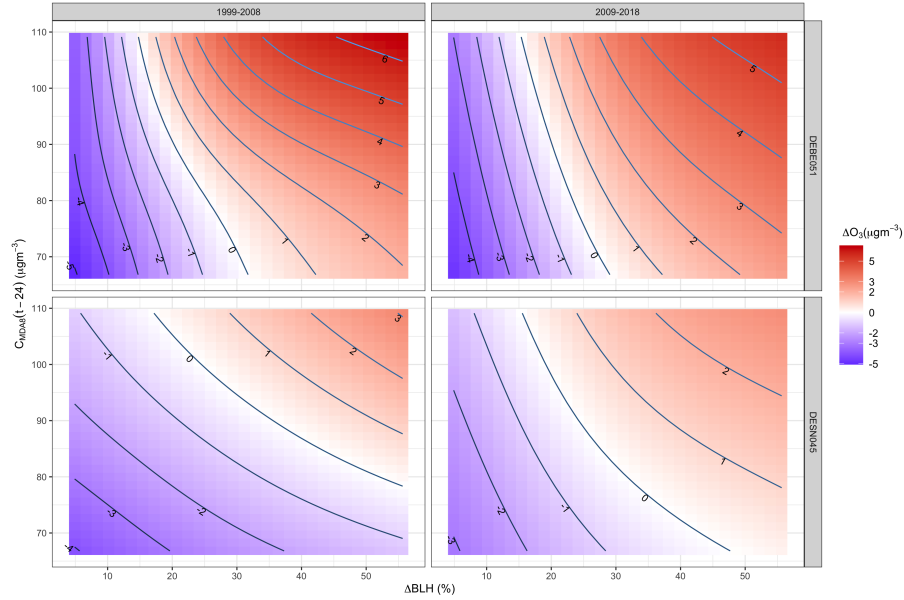


Figure S11. Contour plot for the interaction $\Delta BLH-CMDA_8(t-24)$ at the suburban stations for the periods 1999-2008 and 2009-2018 at the suburban stations.

References

- Akaike, H. 1974. “A New Look at the Statistical Model Identification.” *IEEE T. Automat. Contr.* AC-19: 716–23.
- Barmapadimos, I., C. Hueglin, J. Keller, S. Henne, and A. S. H. Prévôt. 2011. “Influence of Meteorology on Pm10 Trends and Variability in Switzerland from 1991 to 2008.” *Atmos. Chem. Phys.* 11: 1813–35. <https://doi.org/10.5194/acp-11-1813-2011>.
- Boleti, E., C. Hueglin, and S. Takahama. 2019. “Trends of Surface Maximum Ozone Concentrations in Switzerland Based on Meteorological Adjustment for the Period 1990–2014.” *Atmospheric Environment* 213: 326–36. <https://doi.org/10.1016/j.atmosenv.2019.05.018>.
- Carslaw, D. C, S. D Beevers, and J. E Tate. 2007. “Modelling and Assessing Trends in Traffic Related Emissions Using a Generalized Additive Modelling Approach.” *Atmospheric Environment* 41: 5289–99. <https://doi.org/10.1016/j.atmosenv.2007.02.032>.
- Hastie, T., and R. Tibshinari. 1990. “Generalized Additive Models.” *Chapman and Hall, London*.
- Jackson, L. S, N. Carslaw, D. C Carslaw, and K. M Emmerson. 2009. “Modelling Trends in Oh Radical Concentrations Using Generalized Additive Models.” *Atmos. Chem. Phys.* 9: 2021–33. <https://doi.org/10.5194/acp-9-2021-2009>.
- McCullagh, P., and J. A Nelder. 1989. “Generalized Linear Models.” *Chapman & Hall/CRC*. 532.
- Pedersen, E. J, D. L Miller, G. L Simpson, and N. Ross. 2019. “Hierarchical Generalized Additive Models in Ecology: An Introduction with Mgecv.” *PeerJ* 7: e6876. <https://doi.org/10.7717/peerj.6876>.
- R Development Core Team. 2018. “R: A Language and Environment for Statistical Computing.” *Vienna: The R Foundation for Statistical Computing*.
- Wood, S. N. 2006. “Low-Rank Scale-Invariant Tensor Product Smooths for Generalized Additive Mixed Models.” *Biometrics* 62(4): 1025–36. <https://doi.org/10.1111/j.1541-0420.2006.00574.x>.
- . 2011. “Fast Stable Restricted Maximum Likelihood and Marginal Likelihood Estimation of Semi-parametric Generalized Linear Models.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 73(1): 3–36. <https://doi.org/10.1111/j.1467-9868.2010.00749.x>.
- . 2017. “P-Splines with Derivative Based Penalties and Tensor Product Smoothing of Unevenly Distributed Data.” *Statistics and Computing* 27(4): 985–89. <https://doi.org/10.1007/s11222-016-9666-x>.