

Lambda-PFLOTRAN 1.0: Workflow for Incorporating Organic Matter Chemistry Informed by Ultra High Resolution Mass Spectrometry into Biogeochemical Modeling

Katherine A. Muller¹, Peishi Jiang¹, Glenn Hammond¹, Tasneem Ahmadullah¹, Hyun-Seob Song², Ravi Kukkadapu¹, Nicholas Ward³, Madison Bowe³, Rosalie K. Chu¹, Qian Zhao¹, Vanessa A. Garayburu-Caruso¹, Alan Roebuck³, Xingyuan Chen¹

¹ Pacific Northwest National Laboratory, Richland, WA 99352, USA

² Department of Biological Systems Engineering, University of Nebraska—Lincoln, Lincoln, Nebraska, USA

³ Pacific Northwest National Laboratory, Sequim WA 98382, USA

Correspondence to: Katherine Muller (katherine.muller@pnnl.gov)

For submission to Geoscientific Model Development

Abstract. Organic matter (OM) composition plays a central role in microbial respiration of dissolved organic matter and subsequent biogeochemical reactions. Here, a direct connection of organic ~~carbon-matter~~ chemistry and thermodynamics to reactive transport simulators has been achieved through the newly developed Lambda-PFLOTRAN workflow tool that succinctly incorporates carbon chemistry data generated from Fourier transform ion cyclotron resonance mass spectrometry (FTICR-MS) into reaction networks to simulate organic matter degradation and the resulting biogeochemistry. Lambda-PFLOTRAN is a python-based workflow, executed through a Jupyter Notebook interface, that digests raw FTICR-MS data, develops a representative reaction network based on substrate-explicit thermodynamic modeling (also termed lambda modeling due to its key thermodynamic parameter λ used therein), and completes a biogeochemical simulation with the open source, reactive flow and transport code PFLOTRAN. The workflow consists of the following five steps: configuration, thermodynamic (lambda) analysis, sensitivity analysis, parameter estimation, and simulation output and visualization. Two test cases are provided to demonstrate the functionality of the Lambda-PFLOTRAN workflow. The first test case uses laboratory incubation data of temporal oxygen depletion to fit lambda parameters (i.e., maximum utilization rate and microbial carrying capacity). A slightly more complex second test case fits multiple lambda formulation and soil organic matter release parameters to temporal greenhouse gas generation measured during a soil incubation. Overall, the Lambda-PFLOTRAN workflow facilitates upscaling by using molecular-scale characterization to inform biogeochemical processes occurring at larger scales.

31 **1 Introduction**

32 Microbial respiration of dissolved organic carbon (DOC) is a main driver of environmental biogeochemical processes.
33 Mechanistic biogeochemical models often rely on lumping organic matter into a few distinct carbon pools (e.g.,
34 dissolved, sorbed, mineral associated or refractory, labile, etc.) (e.g., Fatichi, et al., 2019, Robertson et al., 2019, Wang
35 et al., 2013) but do not fully consider the properties of the organic matter (OM) compounds individually. Pooled
36 carbon approaches have benefits, such as assigning variable levels of bioavailability, however, this approach does not
37 capture the complex temporal dynamics of respiration driven by OM composition, as aerobic respiration rates have
38 been linked to organic carbon concentration, thermodynamics of the OM (Stegen et al., 2018, Garayburu-Caruso et
39 al., 2020), as well as the diversity of OM compounds present (Lehmann et al. 2020, Stegen et al., 2022). Such findings
40 highlight the importance of incorporating individual OM chemistry into biogeochemical modeling to capture, and
41 ultimately predict, system behavior more accurately.

42 There are many advanced instrumentation techniques capable of detecting and identifying individual OM formulae
43 that comprise a bulk OM sample (e.g., GC-MS, HPLC-MS, Fourier transform ion cyclotron resonance mass
44 spectrometry [FTICR-MS], etc.). For instance, FTICR-MS is a powerful, high-resolution, method that identifies
45 molecular formulae for individual organic compounds. In any given environmental sample, FTICR-MS (or other ultra
46 high-resolution methods) will typically resolve thousands of discrete OM molecular formulae, each with a unique
47 mass and elemental composition (Cooper et al., 2020, Bahureksa et al., 2021). Unfortunately However, untargeted
48 analytical techniques like FTICR-MS are only able to determine if a compound is present and cannot quantify the total
49 concentration associated with each organic matter molecule. Still, such techniques do provide immense amounts of
50 characterization data encompassing a deeper analytical window than measuring a small number of individual
51 biomarkers quantitatively (e.g., Ward et al., 2013). However, the ability to Utilizinge such high-resolution molecular
52 data in reactive transport modeling frameworks affords new opportunity to advance carbon cycling in terrestrial,
53 riverine and coastal systems despite of ~~has remained a challenge and is typically not considered~~ various theoretical and
54 computational challenges.

55 Substrate-explicit thermodynamic modeling (SXTM) provides an avenue for incorporating individual OM reactivity
56 based on thermodynamics (Song et al., 2020) into reactive transport models. The SXTM procedure takes the individual
57 chemical formula derived from FTICR-MS (or another high-resolution technique) and uses its thermodynamic
58 properties to generate an oxidation reaction for each molecular formula present in a sample. The corresponding
59 reaction stoichiometry is then determined by considering catabolic, anabolic, and metabolic reactions and balancing
60 energy for the overall metabolic reaction, allowing for the development of an aerobic respiration expression for each
61 OM formula.

62 Still, the sheer number of compounds identified in each sample proves difficult for model integration. Typically,
63 reactive transport simulators consider only a small number of primary species in their reaction networks, and most
64 could not support modeling each of the thousands of organic matter molecules individually. Here, the developed
65 Lambda-PFLOTRAN workflow addresses this challenge through grouping, or binning, similar compounds based on

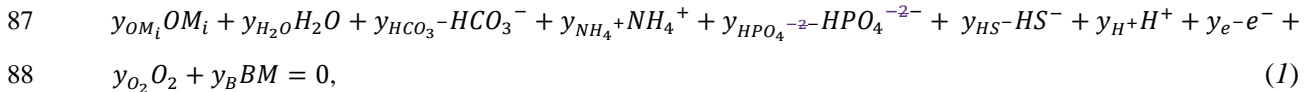
66 their thermodynamic properties, allowing for the number of species considered within the reaction network to be
67 reduced, and thus decreasing the required computational resources.

68 Lambda-PFLOTRAN is a python-based workflow that digests raw FTICR-MS data, develops a representative reaction
69 network based on substrate-explicit thermodynamic modeling (Song et al., 2020), and completes a biogeochemical
70 simulation with the open source, parallel reactive flow and transport code, PFLOTRAN (Hammond et al., 2014).
71 PFLOTRAN is developed under an open source, GNU LGPL license. The term ‘lambda’ is used here because λ is a
72 key parameter in the SXTM, which quantifies thermodynamic favorability of aerobic respiration of OM. The
73 connection between the unique reaction network developed for each FTICR-MS sample hinges on the use of
74 PFLOTRAN’s reaction sandbox capability (Hammond, 2022). The reaction sandbox gives the ability to define
75 additional custom, kinetic reactions beyond standard formulations (e.g., mineral precipitation-dissolution, Michaelis-
76 Menten, etc.). The Lambda-PFLOTRAN workflow enables upscaling by using molecular-scale information to inform
77 larger scale biogeochemical processes occurring throughout a watershed which can be simulated with PFLOTRAN.
78 Herein we describe the Lambda-PFLOTRAN workflow process including the governing expressions, workflow steps,
79 data requirements, as well as the associated assumptions and limitations. Two illustrative test cases are also included
80 to demonstrate the ~~use of the workflow to utilize, parametrize, and model real datasets.~~

81 **2 Methods**

82 **2.1 Conceptual Model**

83 Respiration modeling herein is based on thermodynamic theory by Desmond-Le Quemener and Bouchez (2014) which
84 was updated for multiple OM formulas by Song et al. (2020). The generalized form of OM molecule is assumed to
85 take the form of $C_aH_bN_cO_dP_eS^z_f$. Each molecular formula then undergoes respiration (i.e., reaction with oxygen) based
86 on the following general reaction expression:



89 This generalized expression is used to describe the oxidation of any OM molecule, i , and has been normalized to one
90 mole of biomass (BM) produced. BM is assumed to have a formula of $CH_{1.8}O_{0.5}N_{0.2}$ (Stephanopoulos et al.,
91 1998; Kleerebezem and Van Loosdrecht, 2010). OM_i represents the OM molecules as informed by FTICR-MS. Each
92 y represents the reaction stoichiometry for that reactant ($y < 0$) or product ($y > 0$). While this expression is specific for
93 cases where oxygen is the electron acceptor, such an expression could be updated for alternative electron acceptors.

94 Substrate-explicit thermodynamic modeling expressions developed from Song et al. (2020) were implemented in a
95 reaction sandbox within PFLOTRAN. The expressions were implemented in a general manner allowing for flexibility
96 in handling variations in FTICR-MS data and several user adjustable analysis configurations.

97 The microbial growth kinetics are described by Eq. (2):

$$98 \quad \mu_i^{kin} = \mu^{max} \exp\left(-\frac{\alpha |y_{OM_i}|}{1000 V_h [OM_i]}\right) \exp\left(-\frac{\alpha |y_{O_2}|}{1000 V_h [O_2]}\right), \quad (2)$$

99 where μ_i^{kin} is the unregulated uptake rate of reaction for OM_i [hr^{-1}], μ^{max} is the maximal microbial growth rate [hr^{-1}], y_{OM_i} is the stoichiometry for OM_i [$mol-OM \cdot mol-biomass^{-1}$], V_h is microbial harvest volume [m^3]. Given the
100 physical interpretation of V_h as the microbial harvest volume, it is assumed here that the value of V_h is the same for
101 both OM_i and O_2 , $[OM_i]$ is the organic matter concentration of OM_i [$mol-OM \cdot L^{-1}$], y_{O_2} is the stoichiometry for O_2
102 for respiration of OM_i [$mol-O_2 \cdot mol-biomass^{-1}$], $[O_2]$ is oxygen concentration [$mol-O_2 \cdot L^{-1}$], α is a microbial unit
103 conversion [$mol-biomass$] and 1000 is the conversion of m^3 to L.

104 Further, using a cybernetic modeling approach (after Song et al., 2018), all the unregulated uptake rates (μ_i^{kin}) are
105 normalized by the sum of unregulated uptake rates across all reactions, i following Eq. (3):
106

$$107 \quad u_i = \frac{\mu_i^{kin}}{\sum_{i=1}^n \mu_i^{kin}} \quad (3)$$

108 where u_i is the fraction of the unregulated rate [-]. The final regulated rate, r_i [hr^{-1}] for each reaction is then computed
109 following Eq. (4):

$$110 \quad r_i = u_i \mu_i^{kin}, \quad (4)$$

111 For implementation within PFLOTRAN, the use of inhibition terms was required to prevent negative concentrations
112 once a reactant is nearly depleted. For a reaction to proceed, all reactant species must be present above a minimum
113 concentration even if the molecules do not explicitly control the respiration rate (i.e., species other than OM and O_2 ,
114 Eq. (2). If a reactant concentration falls below a threshold concentration, the respiration rate is inhibited. Reactant
115 inhibition is computed by Eq. 5 (Kinzelbach et al., 1991) for reactant species j :

$$116 \quad I_j = 0.5 + \frac{\arctan([C_j] - C_{thj})f}{\pi}, \quad (5)$$

117 where C_{thj} is the threshold concentration [M], f is the threshold scaling factor [-]. The default C_{thj} is 10^{-20} M.

118 The reaction rates are also inhibited by the microbial carrying capacity of the system, I_{cc} , as follows in Eq. (6):

$$119 \quad I_{cc} = 1 - \frac{[BM]}{CC} \quad (6)$$

120 where [BM] is the biomass concentration [$mol-BM \cdot L^{-1}$], CC is the biomass carrying capacity [$mol-BM \cdot L^{-1}$]. I_{cc} has a
121 non-negativity constraint, so if $[BM] > CC$, then $I_{cc} = 0$.

122 These inhibition factors are applied to the overall rate expression as shown in Eq. (7).

$$r_{i,inhibited} = r_i I_{CC} \prod I_j \quad \forall y_{i,j} < 0, \quad (7)$$

The overall individual species rates, $d[C_j]/dt$, [mol-species·L⁻¹·hr⁻¹] are then computed as follows with Eq. (8):

$$\frac{dC_j}{dt} = (\sum_{i=1}^n y_{i,j} r_{i,inhibited}) [BM], \quad (8)$$

where j is the species index. The total number of species includes 7 general species (i.e., HCO₃⁻, NH₄⁺, HPO₄⁻, HS⁻, H⁺, O₂, BM (i.e., Eq (1)) and the OM species considered (i.e., typically 10). i is the reaction index, n is total number of reactions as based on the total number of OM species (typically, with this workflow $n = 10$). $y_{i,j}$ is the coefficient for species j in reaction i .

The expression for biomass is also modified to account for biomass decay (note all biomass stoichiometries are 1 by definition):

$$\frac{dBM}{dt} = (\sum_{i=1}^n y_{i,j} r_{i,inhibited}) [BM] - k_{deg} [BM], \quad (9)$$

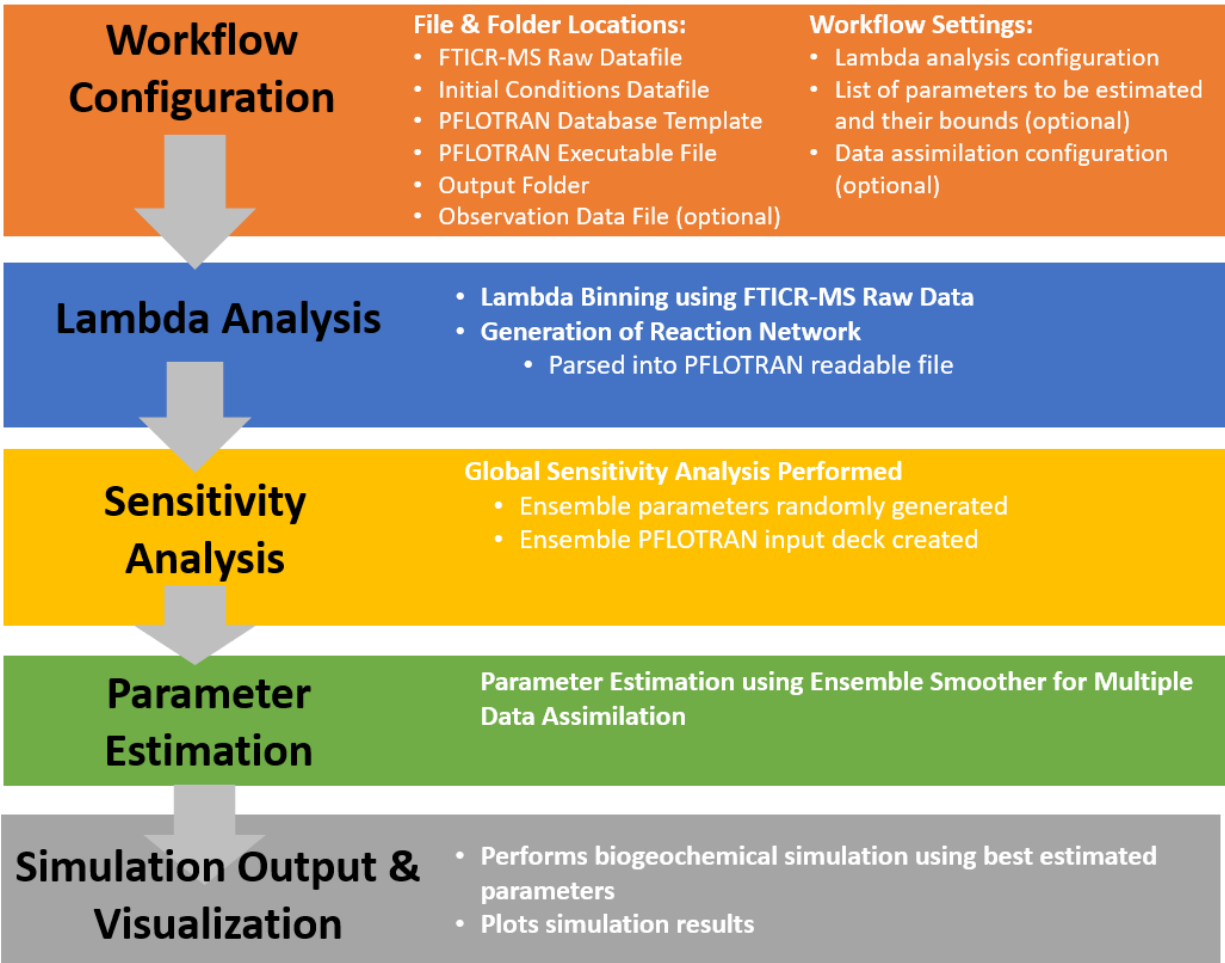
where k_{deg} is the biomass decay rate [hr⁻¹].

2.2 Lambda Analysis and Binning

To reduce the number of organic compounds considered in the simulation, OM molecules are grouped, or binned, based on their λ value computed by Eq. (10):

$$\lambda = \frac{\Delta G_{r,anabolic} + \Delta G_{r,dissipation}}{(-\Delta G_{r,catabolic})}, \quad (10)$$

where ΔG are the Gibbs energies for the anabolic and catabolic reactions and the associated dissipation energy, respectively. The value of λ is indicative of how many times the catabolic reaction needs to be completed to provide the energy required to synthesis one mole of biomass. Lower λ values suggest higher thermodynamic favorability of OM respiration. Using the chemical formula determined for each OM molecule, the energy balance equations are solved providing the overall reaction stoichiometry Eq. (1) and the λ is calculated. Using the λ value for each molecule, the cumulative probability distribution for the sample is produced (Figure 2).

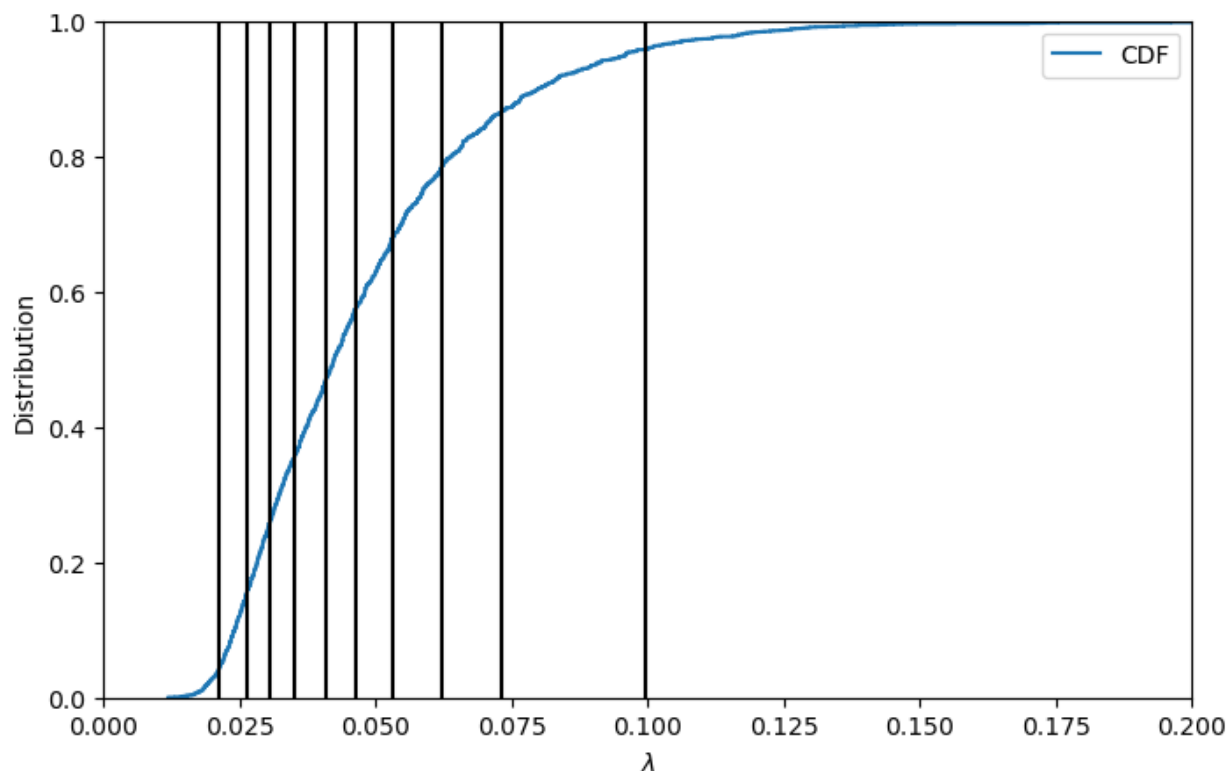


146

147 **Figure 1:** Flow Chart of the Lambda-PFLOTRAN Workflow.

148

149 It is this conversion from individual compounds to a distribution that is critical for reducing the entire sample down
 150 to a representative set of expressions. The λ bins are then formed by splitting the cumulative probability distribution
 151 into equally weighted sections as which to define the overall sample by. The illustrative example shown in Fig. 2
 152 demonstrates the sample distribution being divided into 10 sections (i.e., in this case each section contains 10% of the
 153 overall sample distribution).



154
 155 **Figure 2:** Lambda binning to convert raw FTICR-MS into a representative reaction network using the cumulative probability
 156 distribution function (CDF) for Test Case 1a. Vertical lines display the average λ value for each of the 10 bins (left to right, λ bin
 157 1 to 10).

158
 159 Each section is used to determine a representative organic matter formula and the associated reaction and
 160 stoichiometry of that λ bin. The group of representative reactions (one per bin) is called the reaction network. A
 161 demonstrative reaction network defined by λ analysis and binning is shown in Table 1.

162
 163 **Table 1:** Reaction Network Developed from Lambda Theory for Test Case 1a

Bin Number	Representative Organic Matter Species Formula	λ	y_{OMC}	$y_{HCO_3^-}$	$y_{NH_4^+}$	$y_{HPO_4^{3-}}$	y_{HS^-}	y_{H^+}	y_{O_2}
1	$C_{31}H_{44}N_{0.33}O_{4.8}P_{0.6}S_{0.3}$	0.021	-0.05	0.64	-0.17	-0.18	0.03	0.02	-1.07
2	$C_{26}H_{39}N_{0.20}O_{7.0}P_{0.6}S_{0.1}$	0.026	-0.07	0.68	-0.10	-0.19	0.04	0.01	-1.06
3	$C_{22}H_{36}N_{0.24}O_{7.5}P_{0.5}S_{0.1}$	0.031	-0.08	0.69	-0.02	-0.18	0.04	0.01	-1.06
4	$C_{20}H_{32}N_{0.28}O_{7.3}P_{0.4}S_{0.1}$	0.035	-0.08	0.72	-0.08	-0.18	0.04	0.01	-1.05
5	$C_{19}H_{29}N_{0.48}O_{7.9}P_{0.3}S_{0.2}$	0.041	-0.09	0.79	-0.17	-0.16	0.03	0.02	-1.04
6	$C_{18}H_{26}N_{0.68}O_{8.1}P_{0.2}S_{0.2}$	0.046	-0.10	0.85	-0.27	-0.13	0.02	0.02	-1.03
7	$C_{17}H_{24}N_{0.69}O_{8.1}P_{0.2}S_{0.2}$	0.053	-0.11	0.90	-0.32	-0.12	0.02	0.02	-1.02
8	$C_{15}H_{20}N_{0.67}O_{7.6}P_{0.2}S_{0.2}$	0.062	-0.13	0.94	-0.42	-0.11	0.02	0.03	-1.00

9	$C_{13}H_{19}N_{1.13}O_{87.4}P_{0.1}S_{0.2}$	0.073	-0.15	1.01	-0.48	-0.03	0.01	0.03	-1.00
10	$C_{10}H_{15}N_{1.56}O_{6.5}P_{0.1}S_{0.2}$	0.100	-0.21	1.17	-0.75	0.12	0.01	0.04	-0.97

164
 165 Currently, the representative OM molecule that defines each bin is computed as the average chemical formula of all
 166 the molecules present in that λ section. The disadvantage of this approach is that unrealistic compounds are defined
 167 as representative molecules instead of realistic molecules. The issue with selecting a single, but real compound, from
 168 within each λ section resides in chemical complexity and variation - for instance some molecules may contain low
 169 levels of phosphorous or sulfur and others may not contain either element in the chemical formula. Thus, requiring
 170 the representative chemical formula to be a real compound present in the sample would create basis which would
 171 propagate through the reaction network and into the resulting biogeochemical simulation results.

172 **2.3 Lambda-PFLOTRAN Workflow**

173 The Lambda-PFLOTRAN workflow digests raw FTICR-MS data, calculates the λ distribution for the sample,
 174 generates the λ bins and corresponding reaction network, and completes a biogeochemical simulation using
 175 PFLOTRAN. Further, we incorporated sensitivity analysis and ensemble data assimilation to enable an in-depth
 176 exploration of the impact of reaction parameters on respiration as well as a straightforward parameter estimation
 177 method to fit model parameters to experimental data.

178 The workflow is implemented through a user-friendly Jupyter notebook interface (Kluyver et al., 2016) where a user
 179 can configure the simulation parameters by adjusting initial concentrations, λ binning configuration, parameter values
 180 and/or ranges, and data assimilation options. Based on the user's data file and the associated parameters, scripts within
 181 the Jupyter notebook write the corresponding PFLOTRAN input files, including OM molecules and aqueous
 182 chemistry. The PFLOTRAN simulations are completed locally through a Docker container making this capability
 183 much more user-friendly and accessible. The progress of the data assimilation tool used for parameter fitting is
 184 illustrated within the Jupyter notebook. The resulting best fit final biogeochemical simulation is output visually with
 185 plots and as a text file (when applicable).

186 The Lambda-PFLOTRAN workflow steps are shown in Figure 1 and described in detail in the following subsections:

187 **2.3.1 Step 1 – Workflow configuration**

188 The first step is to set up the workflow configuration for a Lambda-PFLTORAN application. This includes specifying
 189 the file and folder locations of the following information: 1) FTICR-MS raw data file (.csv), 2) initial species
 190 concentrations file (.csv) that includes starting molar concentrations for HCO_3^- , NH_4^+ , HPO_4^{2-} , HS^- , H^+ , O_2 (aq), BM
 191 and total organic carbon (TOC), 3) PFLOTRAN database template file, 4) PFLOTRAN executable file, 5) workflow
 192 output folder, and if completing parameter estimation, (6) the data observation file (.csv), if applicable.

193 The user is also asked to configure workflow settings related to: (1) the lambda analysis configuration, including
 194 number of λ bins and method to define the λ bins (i.e., cumulative vs uniform); (2) the respiration modeling parameter

195 setup, including the list of the parameters to be estimated and their associated upper and lower bounds and (3) the data
196 assimilation configuration (see below).

197 **2.3.2 Step 2 – Organic Matter Chemistry using Lambda Analysis**

198 With only an input of FTICR-MS data, the workflow first performs the lambda analysis (Section 2.2) to group OM
199 molecules into various λ bins based on each compound's thermodynamics (Figure 2) and produce the corresponding
200 reaction network for respiration (Table 1). The default number of λ bins is 10, although this can be adjusted in the
201 workflow configuration by the user, if desired. The generated reaction network is then automatically parsed by the
202 workflow into a text file that can be read by PFLOTRAN.

203 **2.3.3 Step 3 – Sensitivity Analysis using Mutual Information**

204 This step performs the global sensitivity analysis on the parameters to be estimated. Ensemble parameters are first
205 generated by randomly sampling from their predefined ranges in the configuration step and saved into an HDF5 file.
206 Then, the workflow generates a PFLOTRAN input deck to conduct ensemble simulations using the ensemble
207 parameters. The generated ensemble model states enables a global sensitivity analysis using mutual information
208 (Cover and Thomas, 2006; Jiang et al, 2022) as follows:

$$209 \quad I(X;Y) = H(Y) - H(Y|X) = \sum_{x=x} \sum_{y=y} p(x,y) \log \left(\frac{p(x,y)}{p(x)p(y)} \right), \quad (11)$$

210 where x and y are the specific values of X and Y , respectively; $H(Y)$ is the Shannon's entropy of Y ; $H(Y|X)$ is the
211 conditional entropy of Y given X ; p is the probability density function. Higher I indicate stronger sensitivity between
212 X and Y . Besides sensitivity analysis, the ensemble parameter/states also serve as the prior information for parameter
213 estimation at the next step.

214 **2.3.4 Step 4 – Parameter Estimation using Ensemble Smoother for Multiple Data Assimilation**

215 The workflow adopts Ensemble Smoother for Multiple Data Assimilation (Emerick and Reynolds, 2013; Jiang et al,
216 2021), abbreviated as ESMDA, for data assimilation in this step. Rooted in ensemble Kalman filter, ESMDA is an
217 iterative data assimilation approach that assimilates the observations on the entire time period for multiple times to
218 reduce the uncertainty of the estimated or posterior parameters. During each iteration of ESMDA, the model
219 parameters are updated based on the following equation:

$$220 \quad m_{k,l}^u = m_{k,l}^f + C_{MD,l}^f (C_{DD,l}^f + \alpha_l C_D)^{-1} \left(d_{obs} + \sqrt{\alpha_l} C_D^{\frac{1}{2}} z_k - d_{k,l}^f \right), \quad k = 1, \dots, N_e \text{ and } l = 1, \dots, L, \quad (12)$$

221 where the subscripts k and l are the indices of the ensemble member and the iteration, respectively; the superscripts u
222 and f are the updated and forecast parameters or states, respectively; N_e is the number of ensemble members; L is the
223 number of iterations; $m_{k,l}^f$ and $m_{k,l}^u$ are the k th ensemble member of the forecast/prior and updated/posterior
224 parameters, respectively, at the l th iteration; d_{obs} is the observation; z_k is the observation noise sampled from

225 independent standard normal distributions for the k th ensemble member; $d_{k,l}^f$ is the k th ensemble member of the
226 predicted observation states by the model using $m_{k,l}^f$; $C_{MD,l}^f$ is the cross-covariance matrix between the prior parameters
227 m_l^f and the predicted observation states d_l^f ; $C_{DD,l}^f$ is the auto-covariance matrix of the predicted observation states d_l^f ;
228 C_D is the auto-covariance matrix of the observation error; and α_l is the inflation coefficient at the l th iteration with the
229 sum of all α_l equal to one.

230
231 Here, the assimilation starts with taking the ensemble model parameters/states in Step 3 and the provided observations,
232 and calculates the posterior parameters using ensemble Kalman filter, updates the prior parameters with the current
233 posterior for the next iteration, and then repeats the whole process for multiple times (typically 3 to 5 iterations, as
234 defined by the user). The final estimated parameters are obtained from the posterior parameter at the last iteration and
235 are updated in the parameter HDF5 file. The parameter estimation is implemented in a way that allows assimilating
236 either a single (e.g., Test Case 1) or multiple observed species simultaneously through a simple change of the inputs.
237 For example, if temporal experimental or field data is available for oxygen, pH, and total carbon, all these data sources
238 could be simultaneously fit to, with only minor adjustments to Jupyter notebook.

239 2.3.5 Step 5 – Simulation Output and Visualization

240 The last step performs the ensemble simulation of the biogeochemical modeling a final time using the estimated
241 parameters in Step 4. Optionally, users can further pick the realization with the best performance. The user has the
242 option to select their preferred goodness of fit metric from the following options as a means for selecting the best
243 performing simulation: R-squared (R^2), Root Mean Squared Error (RMSE), Modified Kling-Gupta Efficiency
244 ([mMKGE](#)), Nash-Sutcliffe Model Efficiently Coefficient (NSE), or Correlation Coefficient (CorC). Based on the
245 selection, the final time series of aqueous chemistry, oxygen consumption, CO_2 production, and lambda binned, and
246 total organic carbon concentrations will be computed and plotted.

247 3 Test Cases

248 3.1 Test Case 1 - Oxygen Depletion Incubation Experiments.

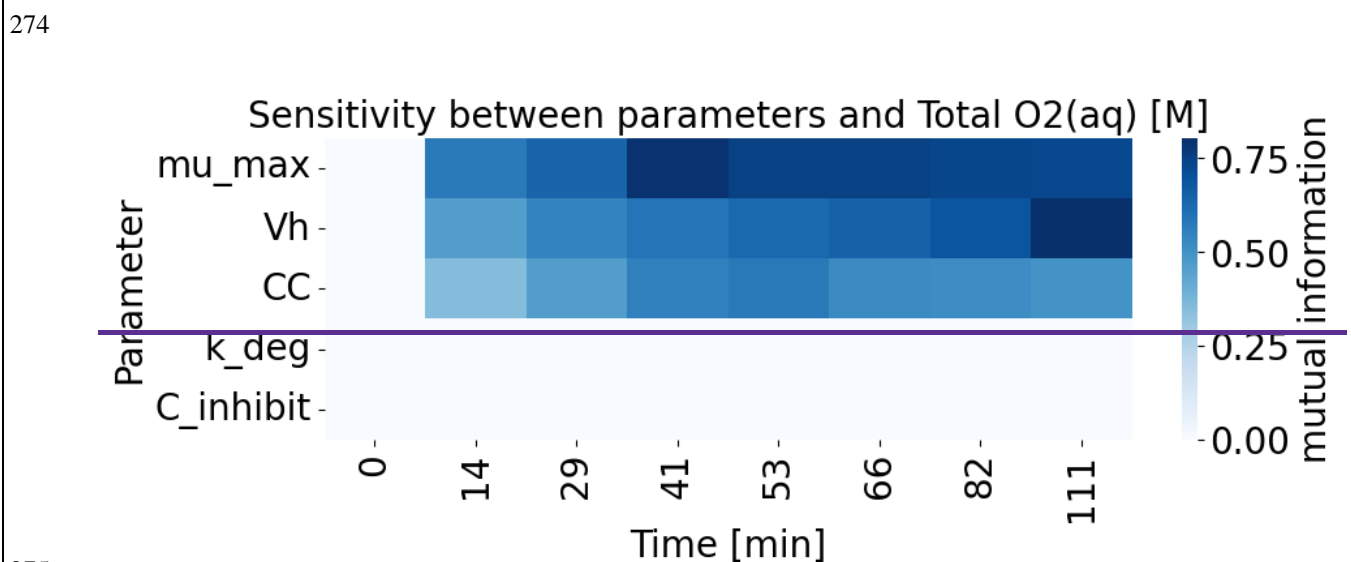
249 In the first illustrative example, the [lambda-pipeline-workflow](#) was used to fit ~~three lambda model parameters~~ (μ_{max} ,
250 V_h , and CC) to laboratory incubation experiments where oxygen levels were measured over two hours in a closed
251 reactor. The incubation experiments were completed as part of the Worldwide Hydrobiogeochemistry Observation
252 Network for Dynamic River Systems (WHONDRS) program (Goldman et al, 2020). For these incubations, sediment
253 was taken from three locations within a stream (i.e., upstream [[Test Case 1a](#)], midstream [[Test Case 1b](#)], and
254 downstream [[Test Case 1c](#)]) in the Yakima River Basin in Washington, USA for subsequent laboratory respiration
255 experiments. FTICR-MS was used to determine the OM chemistry from each sediment sample, resulting in variable
256 formulae being identified in each sample. Formula assignments for all the samples included herein were completed
257 using [formultitude](#) ~~formularity~~ (Tolic et al., 2017). Total dissolved organic carbon concentration paired with the
258 FTICR-MS sample and biomass measurements taken at the start of each experiment were used as the initial

259 concentrations for each of the simulations. Due to the absence of quantitative data related to how the total carbon mass
 260 is distributed between various the OM compounds, the total carbon concentration (on a per-C basis) was assumed to
 261 be split equally between each of the λ bins. The total organic carbon concentration was distributed into each λ bin
 262 using Eq. (13). While this assumption results in equal distribution of carbon between the bins, consequently, it assigns
 263 different initial species concentrations due to varying carbon concentrations between the molecules.

$$264 \quad [C_{\lambda bin}]_0 = \frac{[TOC]}{n_{\lambda bin} n_{C_{\lambda bin}}} \quad (13)$$

265 Where: $[C_{\lambda bin}]_0$ is the initial species concentration in each λ bin [$\text{mol} \cdot \text{L}^{-1}$]; TOC is the total organic carbon measured
 266 [$\text{mol-carbon} \cdot \text{L}^{-1}$]; $n_{\lambda bin}$ is the number of λ bins [-]; and $n_{C_{\lambda bin}}$ is the number of carbon molecules in the assumed
 267 formula for the λ bins [$\text{mol-carbon} \cdot \text{mol-molecule}^{-1}$].

268 Using the Lambda-PFLOTRAN workflow, the FTICR-MS data from each laboratory experiment was digested into
 269 the corresponding λ bins to create the individual reaction network. The Jupyter Notebook for this example is
 270 “Test_Case1-WHONDERS.ipynb” and is available at <https://doi.org/10.15485/2281403>. The resulting λ binning and
 271 associated reaction network for Test Case 1a are shown in Figure 2 and Table 1. Test cases 1b and 1c are in the
 272 Supporting Information (Fig. S1–S2 and Tables S2–S3). The calculated parameter sensitivity is shown in Figure 3,
 273 which indicates the results highly sensitive to all three parameters, in particular μ_{max} and V_h , more so than and CC .

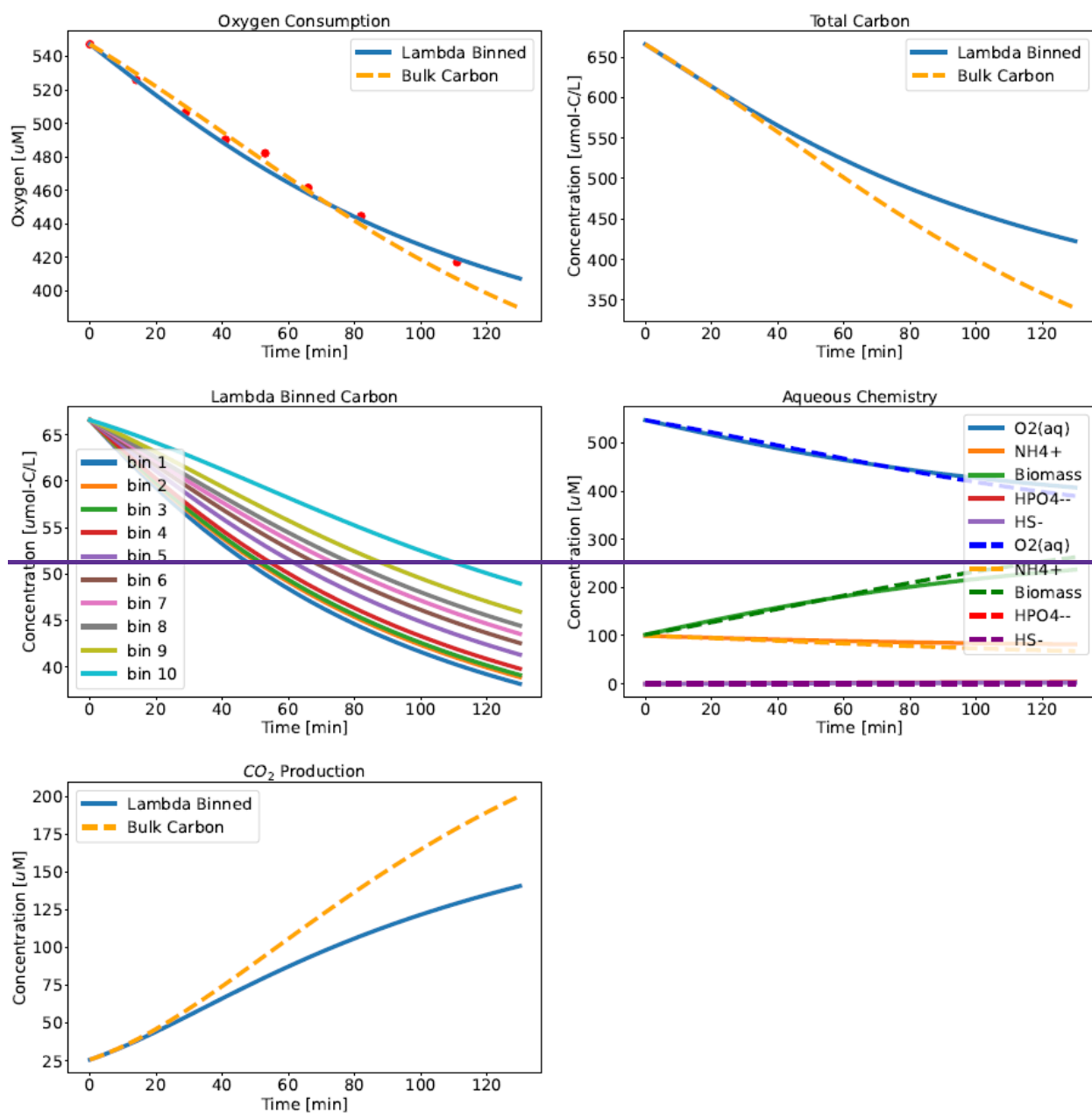


275

276 **Figure 3:** Illustrative Example of Sensitivity Analysis Output during Parameter Estimation. Example shown here provides the
 277 sensitivity of three fitted parameters (μ_{max} , V_h , and CC) on temporal aqueous O_2 concentrations as a function of time.

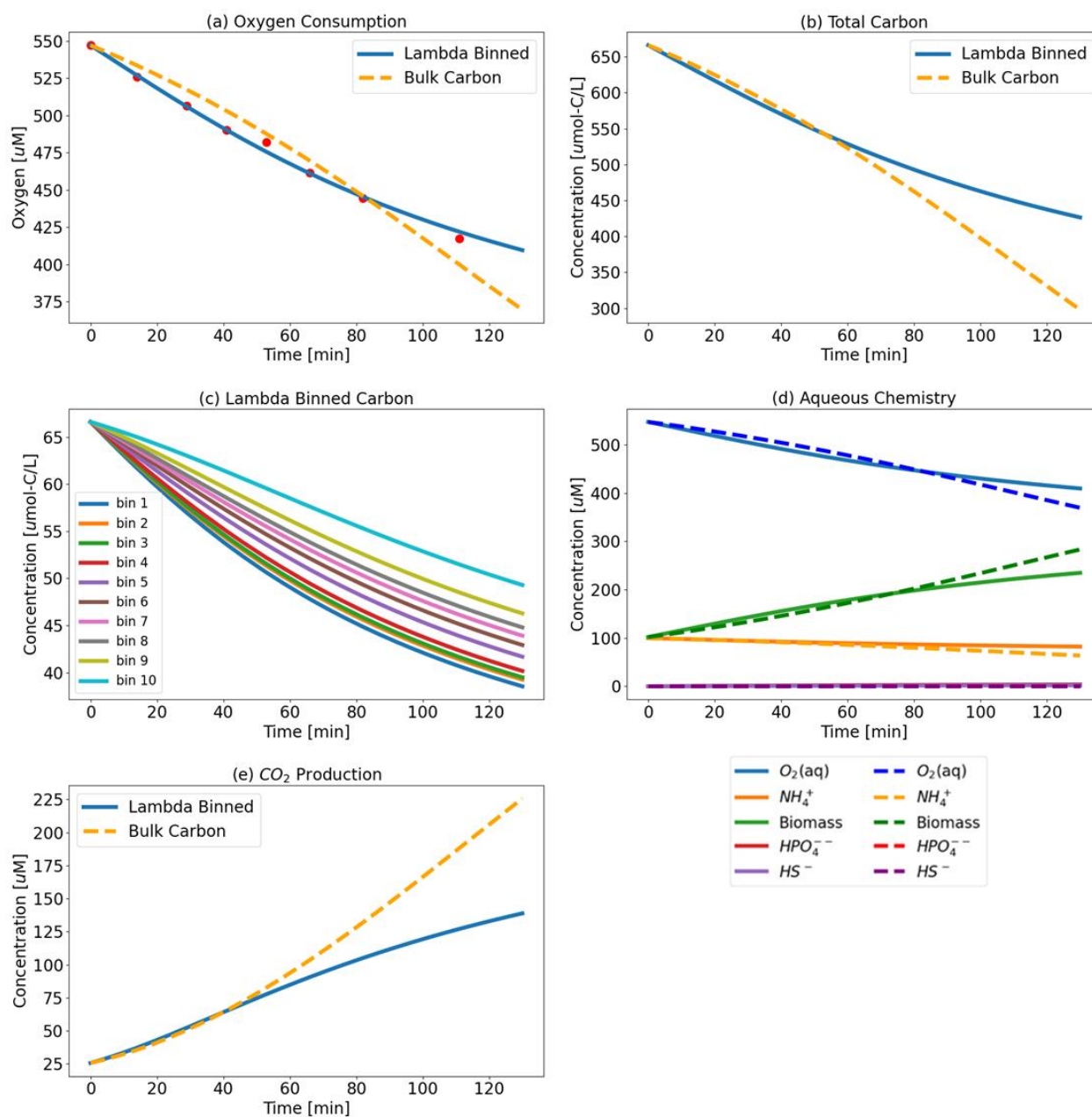
278 μ_{max} Lambda expression parameters were fit to the provided experimental oxygen data. The final lambda binned
 279 fit, along with and the final fit to the experimental data and corresponding carbon consumption (individual and total)
 280 and aqueous chemistry is displayed in Figure 4-3 (and in the supporting information Fig. S2-S1 and S3-S2 for Test

281 Cases 1b (midstream) and 1c (downstream), respectively). The workflow was also run where μ_{\max} was fit again, but
 282 this time assuming a generic OM form of CH_2O . This, allowing To evaluate the allows for comparison between using
 283 information for use of lambda binned OM obtained from FTICR-MS (Figure 43). The workflow was also run for a
 284 baseline case where μ_{\max} was fit again, but this time assuming a generic bulk OM form of CH_2O for comparison. Fitted
 285 μ_{\max} values for the lambda binned model is 0.25 min^{-1} ($R^2 = 0.99$) and fitted μ_{\max} to the bulk OM CH_2O model is 0.032
 286 min^{-1} ($R^2 = 0.96$). V_h and CC are fixed at assumed values of 10 m^3 and 1 M , respectively in both simulations. In this
 287 case, the same set of lambda parameters were fit to the oxygen consumption experimental data, which also resulted in
 288 successful fit ($R^2 = 0.990$ for lambda binned model; $R^2 = 0.987$ for bulk model).



289

290 **Figure 4:** Test Case 1a Results—Oxygen Consumption (top left) where Lambda-PFLOTTRAN workflow was used to fit (blue line)
 291 to experimental respiration data (red dots) and the corresponding Total Carbon Consumption (top right); Individual Organic Matter
 292 Consumption by λ bin (middle left); corresponding biogeochemistry including O_2 (aq) (blue); Biomass (green); NH_4^+ (orange);
 293 HS^- (purple); and HPO_4^{2-} (red) (middle left); and CO_2 production (bottom left) for the upstream incubation. The dashed orange
 294 lines in the top two figures show simulation results assuming a generic OM species of CH_2O for comparison. Fitted values for the
 295 lambda binned model are $\mu_{max}=0.25 \text{ min}^{-1}$, $V_h=9.7 \text{ m}^3$, and $CC=0.49 \text{ M}$ ($R^2=0.990$), and fitted bulk-OM CH_2O model value
 296 are $\mu_{max}=0.05 \text{ min}^{-1}$, $V_h=3.3 \text{ m}^3$, and $CC=0.58 \text{ M}$ ($R^2=0.987$).



297
 298 **Figure 3:** Test Case 1a Results – (a) Oxygen Consumption where Lambda-PFLOTTRAN workflow was used to fit (blue line) to
 299 experimental respiration data (red dots) and (b) Total Carbon Consumption; (c) Individual Organic Matter Consumption by λ bin;
 300 and (d) biogeochemistry including O_2 (aq) (blue); Biomass (green); NH_4^+ (orange); HS^- (purple); and HPO_4^{2-} (red); and (e) CO_2

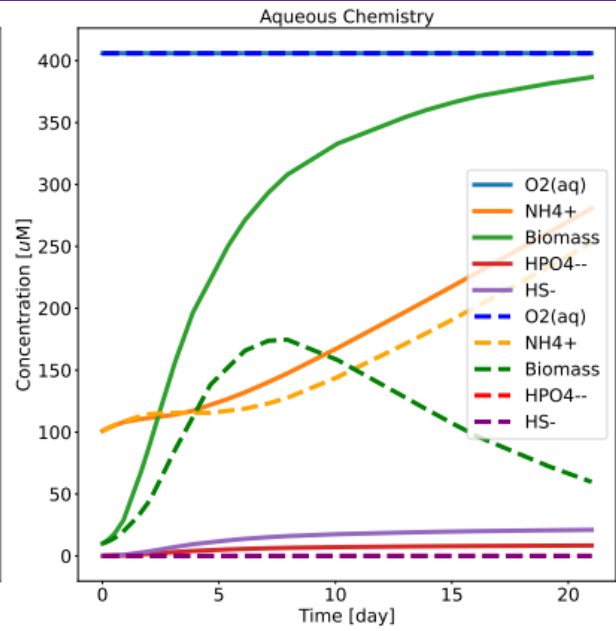
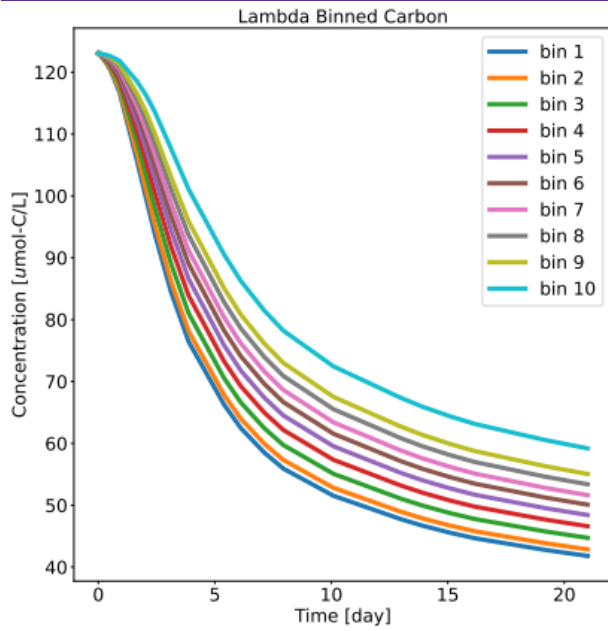
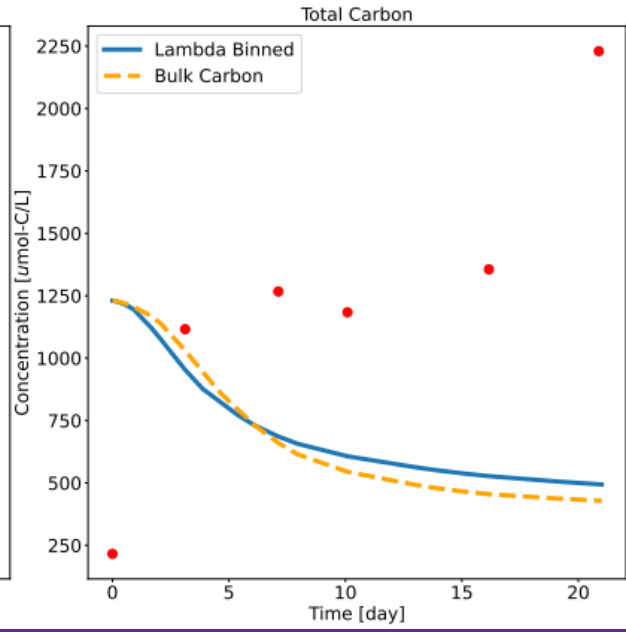
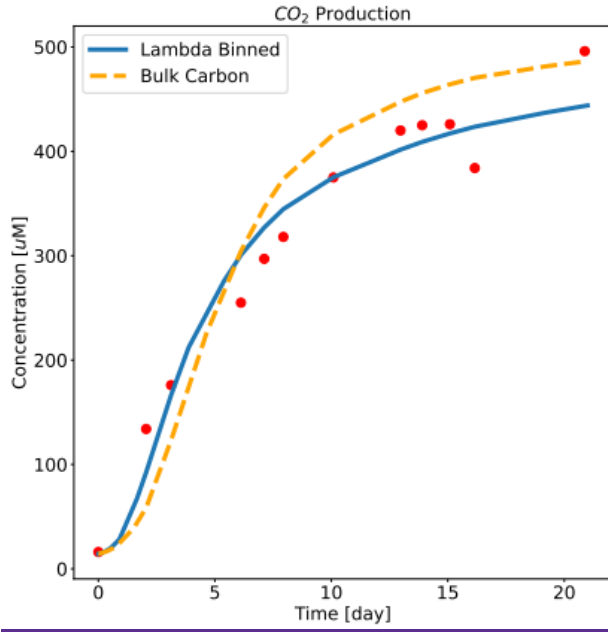
301 [production for the upstream incubation. The dashed orange lines \(in a, b and e\) show simulation results assuming a generic OM](#)
302 [species of CH₂O for comparison.](#)

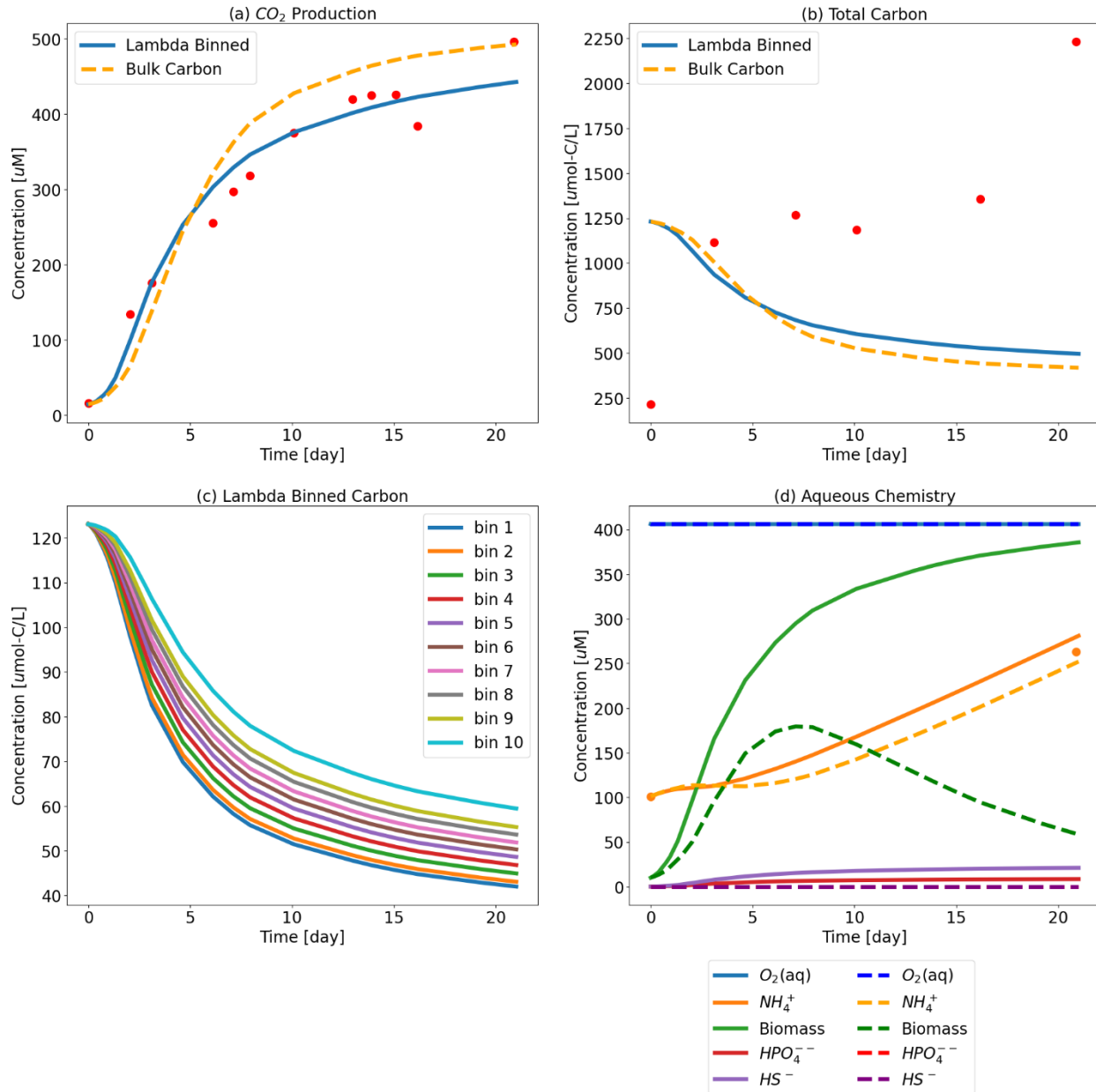
303 However, even over the short time frame of this simulation ([i.e., only 120 minutes](#)), the difference between assuming
304 the generic CH₂O and using the more detailed organic matter chemistry resulted in different predictions of total carbon
305 and CO₂ generation. The bulk OM model predicts more carbon consumption and greater CO₂ production than the
306 ~~binned~~-lambda [binned](#) model. The bulk OM model estimates that ~~65~~[50](#)% of the initial total carbon is consumed over
307 the first 120 mins, whereas the lambda binned model predicts ~~54~~[34](#)% consumption. Similarly, the bulk OM model
308 estimates approximately ~~41~~[35](#)% more CO₂ generation as compared the lambda binned model. The effects on aqueous
309 chemistry over this short duration are more muted, albeit still present.

310 **3.2 Test Case 2 - Respiration Incubation Experiments.**

311 Test Case 2 uses soil respiration incubation data from Ward et al. (2023) aimed at investigating the influence of soil
312 type, oxygen condition (aerobic vs. anaerobic), and seawater exposure (fresh vs. saline) on respiration extent and rate.
313 For these experiments, temporal measurements were collected for CO₂ generation, dissolved organic carbon (DOC),
314 organic matter formulas via FTICR-MS and other bulk aqueous chemistry (i.e., pH, NH₄⁺, and other metals and ions)
315 creating a rich dataset for calibration of system specific lambda model parameters. These incubations were setup by
316 adding dry soil to the reactor and then adding water (resulting in a soil:water ratio ranging from 1:11 to 1:16). The soil
317 and water were shaken vigorously for five minutes, and then sampled for the initial time point prior to officially
318 starting the incubation. For the aerobic experiments, the reactor headspace was cycled every 24 hours to measure CO₂
319 generated but also to ensure the system was kept aerobic; this was only performed five days per week, with no
320 measurements taken on the weekend due to logistical constraints. Upon experiment completion, the increase in DOC
321 concentrations indicated organic carbon was being kinetically released from the soil into the aqueous phase over the
322 course of the 21-day experiment. Similarly, measured NH₄⁺ concentrations also increased during the experiment. To
323 address this [within our reactive transport model](#), a source of nitrogen was assumed to be released from the soil as well
324 ($N_{release}$). Both carbon and nitrogen release are included in this example and are assumed to follow a zero-order constant
325 release rate. Any organic carbon released from the soil was fractionated into each λ bin on the same per-carbon basis
326 assumed for the initial total organic carbon. This was implemented through a dependent function that calculated the
327 release of carbon into each λ bin based on a fitted single bulk $k_{release}$ rate. Mathematically in PFLOTRAN the constant
328 oxygen conditions were implemented through a gas-liquid partitioning expression with a fast exchange term. [These](#)
329 [three additional processes were added to describe the experimental conditions of Test Case 2 more accurately \(i.e.,](#)
330 [release of carbon, nitrogen and sustained aerobic conditions\)](#); however, a PFLOTRAN input deck can be expanded
331 [and customized to include a host of additional processes and full geochemistry for a specific system of interest. For](#)
332 [instance, aqueous complexation, mineral dissolution and precipitation, sorption, and redox reactions can be added, all](#)
333 [of which can influence the resultant pH and carbon, nitrogen, and other nutrient dynamics.](#)

334 The workflow was used to fit μ_{max} , V_h , CC, k_{deg} , as well as $k_{release}$, to the temporal CO₂ generation for a single aerobic
335 soil incubation (Figure 5). The Jupyter Notebook for this example is “Test_Case2-Colloids.ipynb”.



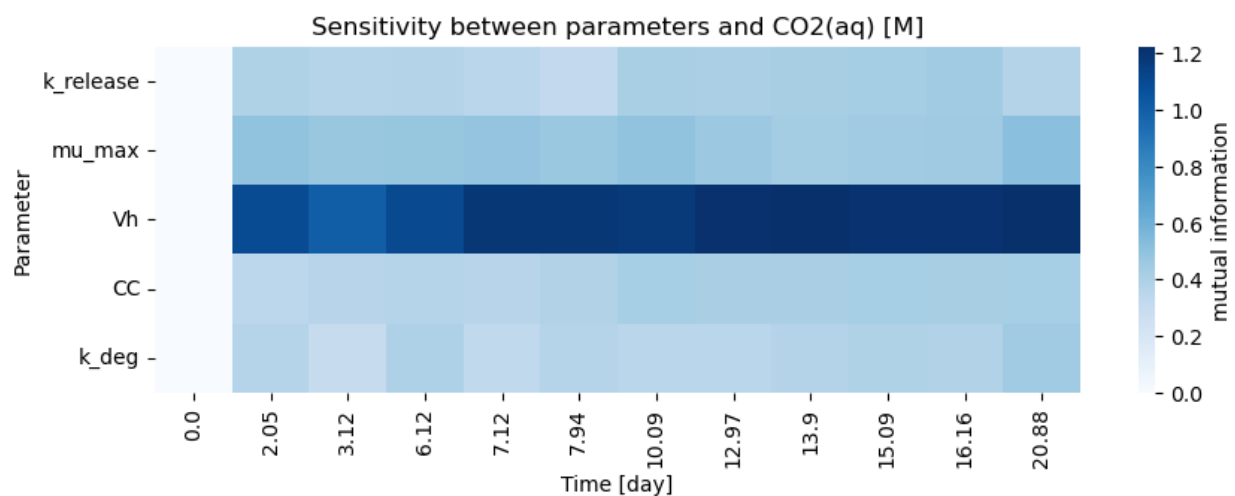


337

338 **Figure 54.** Test Case 2 Results – (a) CO₂ production [top-left](#) where Lambda-PFLOTRAN workflow was used to fit (blue line) to
 339 experimental respiration data (red dots) and (b) the corresponding Total Organic Carbon [top-right](#); (c) Individual Organic Matter
 340 Consumption by λ bin [bottom-left](#) and (d) the corresponding biogeochemistry including O₂ (aq) (blue); Biomass (green); NH₄⁺
 341 (orange); HS⁻ (purple); and HPO₄⁻ (red) [bottom-right](#). Dots indicate experimental data. The dashed orange lines in the top two
 342 figures show simulation results assuming a generic OM species of CH₂O for comparison. Fitted parameters for lambda binned
 343 model were $k_{\text{release}} = 5.5 \times 10^{-12} \text{ day}^{-1}$; $\mu_{\text{max}} = 37.6 \text{ day}^{-1}$, $V_h = 5.0 \text{ m}^3$, $CC = 0.12 \text{ M}$, and $k_{\text{deg}} = 1 \times 10^{-3} \text{ day}^{-1}$ ($R^2 = 0.953$) and fitted
 344 bulk OM CH₂O model values were $k_{\text{release}} = 2.0 \times 10^{-12} \text{ day}^{-1}$; $\mu_{\text{max}} = 47 \text{ day}^{-1}$, $V_h = 1.0 \text{ m}^3$, $CC = 0.77 \text{ M}$, and $k_{\text{deg}} = 0.15 \text{ day}^{-1}$ ($R^2 =$
 345 0.909).

346 [The best fit results indicate a superior fit using the lambda binned OM over the bulk OM model and in fact, the bulk](#)
 347 [model is unable to successfully capture the temporal evolution of the CO₂. It should be noted that both model fits are](#)

348 highly sensitive to the allowable parameter space as user defined by the lower and upper parameter bounds. For the
 349 purposes for showcasing the workflow, five parameters were estimated in this test case example, and as a result the
 350 models are likely over parametrized given the amount of data available. Parameter sensitivity over the course of
 351 simulation time is shown in Figure 5 and suggests this system is highly sensitive to V_h . It should be noted that both
 352 these model fits are also highly sensitive to the allowable parameter space as user defined by the lower and upper
 353 parameter bounds.



354
 355 **Figure 5.** Test Case 2 - Sensitivity Analysis Output during Parameter Estimation. The sensitivity of five fitted parameters (k_{release} ,
 356 μ_{max} , V_h , CC, and k_{deg}) on temporal aqueous CO₂ concentrations as a function of time.

357 Any additional experimental data, either collected during incubations or through independent experiments (e.g.,
 358 carbon release from the soil in an abiotic system), would be expected to help constraint the model and improve
 359 parameterization. Additionally, it is unclear why the model is unable to capture the total organic carbon behavior in
 360 Test Case 2. One potential explanation is that some of the released organic carbon may not be fully bioavailable and
 361 thus the model may be compensating for this by artificially reducing the concentration of OM available for respiration.

362 **4 Variability and Impact of Organic Matter Speciation**

363 The variability in OM speciation was briefly assessed by comparing FTICR-MS data from Test Cases 1 and 2. Each
 364 identified OM species was classified into one of nine compound classes. For Test Case 1, the average of the three Test
 365 Case 1 samples (1a - upstream, 1b - midstream, and 1c - downstream) was computed. The predominant classes were
 366 proteins ($34 \pm 1\%$), lignin ($26 \pm 1\%$), and lipids ($13 \pm 2\%$), with the errors representing the standard deviation among
 367 the Test Case 1a-c samples. The low standard deviation suggests consistent reproducibility in OM speciation for
 368 samples taken from nearby locations. In contrast, OM in Test Case 2 was primarily composed of lignin (37.4%) and
 369 concentrated hydrocarbons (32%). The full distribution of compound classes is presented in Figure 5.

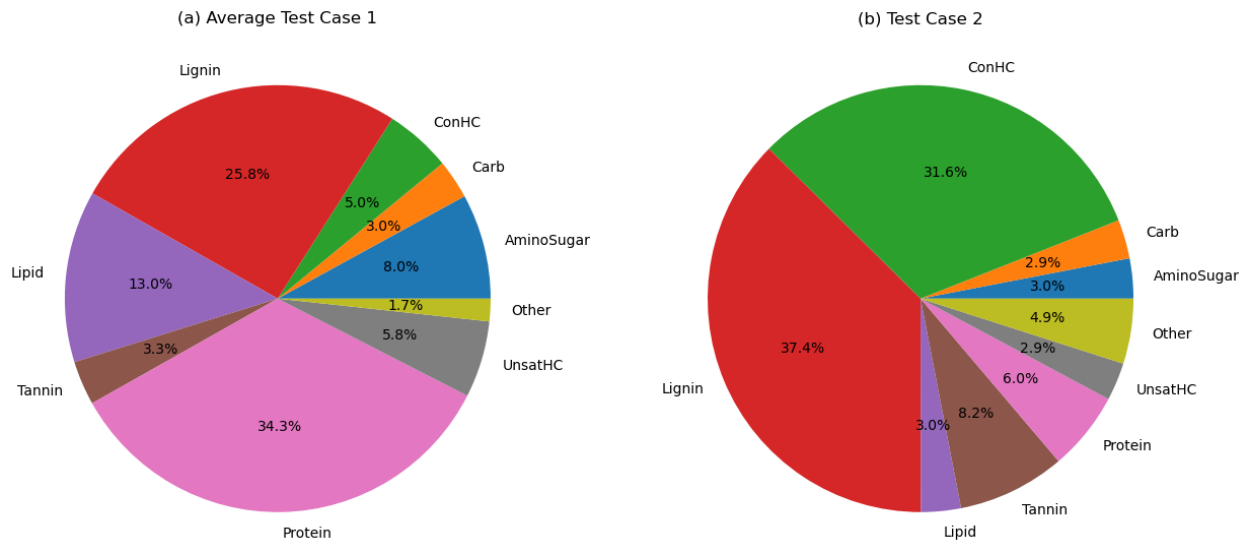
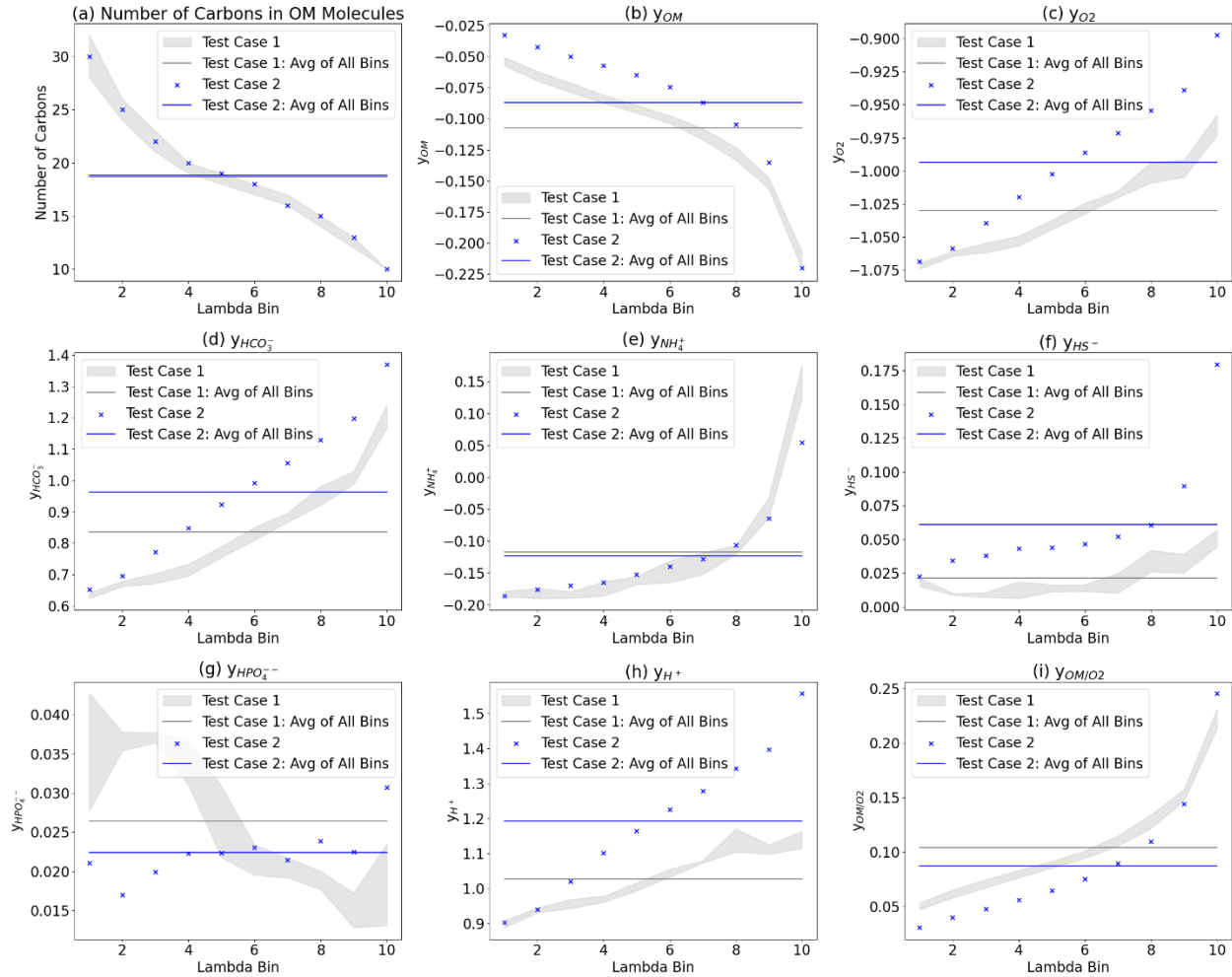


Figure 5. Distribution of Organic Matter Compound Classes: (a) Test Case 1 and (b) Test Case 2.

Note: Test Case 1 is the average of Test Case samples 1a-c. ConHC = Condensed Hydrocarbon; UnsathC = Unsaturated Hydrocarbon

The influence of the sample OM speciation on the λ binned reaction networks was also assessed. Figure 6 illustrates the impact of OM speciation on the corresponding λ binned reaction networks, with three key observations. First, the variability in OM speciation between different samples is evident when comparing Test Case 1 and Test Case 2. To enhance visual clarity, the range of Test Case 1 samples (1a-c) is depicted as a grey shaded region, showing the spread between the minimum and maximum values of the three samples. For Test Case 2, data from the single FTICR-MS sample is represented by blue dots. Test Case 1 and 2 have distinct λ derived reaction networks as indicated by the little overlap between the grey region and the blue dots in Figures 6b-i.



382

383 **Figure 6.** Comparison of Lambda-Binned Reaction Network Parameters: (a) number of carbons in the OM; stoichiometric
 384 coefficient, y , for (b) OM, (c) O_2 , (d) HCO_3^- , (e) NH_4^+ , (f) HS^- , (g) HPO_4^- , (h) H^+ ; and (i) ratio of OM/ O_2 coefficients for Test
 385 Case 1a-c (grey dots); the average of all λ bins for Test Case 1 (grey line); Test Case 2 (blue x); and the average of all λ bins for
 386 Test Case 2 (blue line). The grey shaded area highlights the range of values for Test Case 1a-c for better visual comparison.

387

388 Second, the λ binning process captures the OM speciation variation within a sample. To illustrate this intrasample
 389 variability, a line representing the average of all λ bins is shown on Figure 6 (grey line for Test Case 1, blue line for
 390 Test Case 2). The difference between the reaction network coefficients (vertical axis) for the λ binning (grey shaded
 391 area and blue dots) and the Test Case average lines highlights the extent of this variability. Finally, although the λ
 392 binning process resulted in a similar number of carbon atoms to OM molecules within each λ bin for both test cases
 393 (Figure 6a), the resulting stoichiometric coefficients in the reaction networks differ significantly (Figures 6b-h). These
 394 stoichiometric differences lead to variations in biogeochemical outcomes, such as OM-to-oxygen utilization ratios
 395 during aerobic respiration (Figure 6i). These differences are due to the additional elements beyond carbon in the OM
 396 molecules (i.e., nitrogen, oxygen, sulfur, hydrogen, and phosphorus).

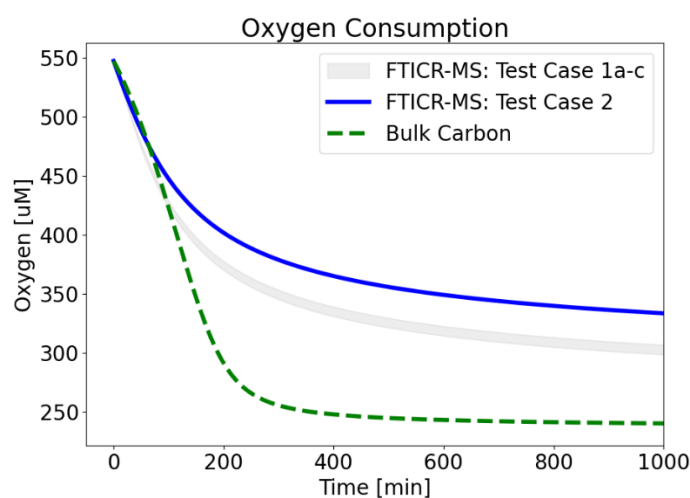
397

398 Another important aspect is the comparison against assuming a generic, bulk OM composition, which does not account
399 for OM speciation as informed by FTICR-MS or similar methods. The reaction network developed for a generic OM
400 molecule of CH₂O is shown in Eq. 14.



402 This reaction network is used in the Lambda-PFLOTRAN workflow for bulk OM simulations.
403 To further assess and isolate the effect of OM speciation, extended forward simulations were performed by only
404 varying FTICR-MS input data (Figure 7). FTICR-MS samples from Test Cases 1a-c and Test Case 2 were tested.
405 These simulations replicate Figure 3 (i.e., Test Case 1a conditions and fitted μ_{max} values) with the expectation of OM
406 speciation, and demonstrate the significant impact of OM chemistry and speciation on overall predicted behavior,
407 especially over longer time periods.

408



409
410 Figure 7. Influence of OM Speciation on Oxygen Consumption. FTICR-MS data from Test Cases 1a-c (grey shaded area), and
411 Test Case 2 (blue line) were used as inputs. Bulk CH₂O OM (green line) was also plotted for reference. Best fit μ_{max} values to Test
412 Case 1a were used (i.e., lambda binned $\mu_{\text{max}} = 0.25 \text{ min}^{-1}$; bulk OM $\mu_{\text{max}} = 0.032 \text{ min}^{-1}$).

413
414 The clear variability in OM speciation, differences between a generic OM reaction network and one informed by
415 FTICR-MS, and the impact of OM chemistry on biogeochemical predictive simulations underscore the importance of
416 incorporating site-specific OM chemistry informed by ultra high resolution characterization into biogeochemical
417 models.

418 **4.5 Conclusions**

419 Overall, Lambda-PFLOTRAN workflow provides an important linkage between molecular scale organic matter
420 characterization and reactive transport simulations. This workflow allows for the influence of organic matter
421 composition to be utilized within simulators to provide a more comprehensive understanding of the system chemistry

422 and behavior, moving beyond the standard assumption of bulk organic matter chemistry and composition. While there
423 are current limitations due to how composition is characterized and quantified, this workflow connecting
424 characterization information to simulations is an important advancement that can be refined as these laboratory
425 techniques improve over time.

426 One of the major limitations surrounding this method, is the lack of understanding of organic matter compound
427 bioavailability, resulting in a large conceptual gap as to how various organic carbon compounds may be utilized by
428 microbes. In the absence of such information, all identified organic matter molecules are assumed to have equal
429 bioavailability within this modeling framework when, in reality, compounds will exhibit varying degrees of
430 bioavailability depending on factors such as associated size fraction, carbon pool, and environmental factors (Schmidt
431 et al., 2011; Ahamed et al., 2023). Until improved understanding is established to discern individual compound
432 bioavailability, this will remain as a limitation.

433 Another limitation of this method resides around the analytical limitations of organic carbon characterization and
434 quantification. For instance, FTICR-MS focuses on water soluble organic matter which may provide a bias in the
435 types of carbon identified by this technique (Tfaily et al., 2017). Additionally, as mentioned previously, FTICR-MS
436 is qualitative, it does not provide structural information and will not differentiate between different isomers that have
437 the same molecular formulas, it is only able to identify molecular formula is present or absent and not the concentration
438 associated with each peak. Here, this has been addressed by assuming equal distribution of total carbon between the
439 formulas within each λ bin on a per-carbon basis. This caveat can be easily updated in the workflow if new analytical
440 advances are made that provide more quantitative information. Some existing approaches could be suitable for this
441 type of modeling such as using quantitative biomarkers that cover major compound classes (Kim and Blair, 2023);
442 but further advances in obtaining both high resolution and quantitative OM characterization would greatly aid in how
443 we understand and model ecosystems.

444

445 **Acknowledgements:**

446 This research was performed under a variety of interdisciplinary projects including the U.S. Department of Energy
447 (DOE) sponsored Office of Science, Office of Biological and Environmental Research (BER), Environmental System
448 Science (ESS) Program, IDEAS-Watersheds, River Corridor Scientific Focus Area (SFA), the Environmental
449 Molecular Sciences Laboratory User Facility sponsored by the Biological and Environmental Research program under
450 Contract No. DE-AC05-76RL01830, and COMPASS-FME, a multi-institutional project supported by DOE-BER as
451 part of the Environmental System Science Program. This study used data from the Worldwide Hydrobiogeochemistry
452 Observation Network for Dynamic River Systems (WHONDRS). This paper describes objective technical results and
453 analysis. The work was performed at the Pacific Northwest National Laboratory (PNNL). [PNNL is operated for DOE
454 by Battelle Memorial Institute under contract DE-AC05-76RL01830. This paper describes objective technical results
455 and analysis.](#) Any subjective views or opinions that might be expressed in the paper do not necessarily represent the
456 views of the U.S. Department of Energy or the United States Government.

457

458 **Code Availability:**

459 The source code, installation requirements, example test case notebooks, and associated data are available in ESS
460 DIVE at <https://doi.org/10.15485/2281403>

461

462 **Author Contribution:**

463 KM: conceptualization, formal analysis, methodology, software, writing- original draft preparation; PJ: methodology,
464 software, writing- original draft preparation; GH: methodology, software, writing-review & editing; TA: data curation,
465 software, writing-review & editing; HS: methodology, writing-review & editing; RK: supervision; NW: supervision,
466 writing-review & editing; MB: investigation; RC: investigation; QZ: investigation; VG: investigation, data curation;
467 AR: investigation; XC: conceptualization, investigation, writing-review & editing

468

469 **Competing Interests:** The authors declare that they have no conflict of interest.

470 **References**

471 Bahureksa, W., Tfaily, M. M., Boiteau, R. M., Young, R. B., Logan, M. N., McKenna, A. M., & Borch, T. (2021).
472 Soil organic matter characterization by Fourier transform ion cyclotron resonance mass spectrometry (FTICR MS): A
473 critical review of sample preparation, analysis, and data interpretation. *Environmental science & technology*, 55(14),
474 9637-9656, <https://doi.org/10.1021/acs.est.1c01135>
475 Cover, T. M., and Thomas, J. A. (2006). Elements of information theory (Wiley series in telecommunications and
476 signal processing). Wiley-Interscience.
477 Emerick, A.A., Reynolds, A.C. (2013). Ensemble smoother with multiple data assimilation. *Comput. Geosci.* 55, 3–
478 15. <https://doi.org/10.1016/j.cageo.2012.03.011>.

479 Fatichi, S., Manzoni, S., Or, D., & Paschalis, A. (2019). A mechanistic model of microbially mediated soil
480 biogeochemical processes: a reality check. *Global Biogeochemical Cycles*, 33(6), 620-648.
481 <https://doi.org/10.1029/2018GB006077>

482 Garayburu-Caruso, V., Stegen, J., Song, H.-S., Renteria, L., Wells, J., Garcia, W., et al. (2020). Carbon limitation
483 leads to thermodynamic regulation of aerobic metabolism. *Environ. Sci. Technol. Lett.* 7, 517–524.
484 <https://doi.org/10.1021/acs.estlett.0c00258>

485 Goldman A E ; Arnon S ; Bar-Zeev E ; Chu R K ; Danczak R E ; Daly R A ; Delgado D ; Fansler S ; Forbes B ;
486 Garayburu-Caruso V A ; Graham E B ; Laan M ; McCall M L ; McKeever S ; Patel K F ; Ren H ; Renteria L ; Resch
487 C T ; Rod K A ; Tfaily M ; Tolic N ; Torgeson J M ; Toyoda J G ; Wells J ; Wrighton K C ; Stegen J C ; WHONDRS
488 Consortium T (2020): WHONDRS Summer 2019 Sampling Campaign: Global River Corridor Sediment FTICR-MS,
489 Dissolved Organic Carbon, Aerobic Respiration, Elemental Composition, Grain Size, Total Nitrogen and Organic
490 Carbon Content, Bacterial Abundance, and Stable Isotopes (v8). River Corridor and Watershed Biogeochemistry SFA,
491 ESS-DIVE repository. Dataset. doi:10.15485/1729719 accessed via [https://data.ess-](https://data.ess-dive.lbl.gov/datasets/doi:10.15485/1729719)
492 [dive.lbl.gov/datasets/doi:10.15485/1729719](https://data.ess-dive.lbl.gov/datasets/doi:10.15485/1729719) on 2023-12-28

493 Hammond, G. E., Lichtner, P. C., & Mills, R. T. (2014). Evaluating the performance of parallel subsurface simulators:
494 An illustrative example with PFLOTRAN. *Water resources research*, 50(1), 208-228.
495 <https://doi.org/10.1002/2012WR013483>

496 Hammond, G.E. (2022) The PFLOTRAN Reaction Sandbox, *Geoscientific Model Development*, 15, 1659-1676,
497 <https://doi.org/10.5194/gmd-15-1659-2022>.

498 Jiang P., Chen, X., Chen, K., Anderson, J., Collins, N., Gharamti, M. (2021) DART-PFLOTRAN: An ensemble-based
499 data assimilation system for estimating subsurface flow and transport model parameters. *Environmental Modelling &*
500 *Software*, Volume 142, <https://doi.org/10.1016/j.envsoft.2021.105074>.

501 Jiang, P., Son, K., Mudunuru, M.K. and Chen, X. (2022). Using mutual information for global sensitivity analysis on
502 watershed modeling. *Water Resources Research*, 58(10), <https://doi.org/10.1029/2022WR032932>

503 Kim, J., Blair, N.E. Biomarker heatmaps: visualization of complex biomarker data to detect storm-induced source
504 changes in fluvial particulate organic carbon. *Earth Sci Inform* 16, 2915–2924 (2023). [https://doi.org/10.1007/s12145-](https://doi.org/10.1007/s12145-023-01039-y)
505 [023-01039-y](https://doi.org/10.1007/s12145-023-01039-y)

506 Kinzelbach, W., Schafer, W., and Herzer, J. (1991). Numerical modeling of natural and enhanced
507 denitrification processes in aquifers. *Water Resources Research*, 27(6):1123–1135.
508 <https://doi.org/10.1029/91WR00474>

509 Kleerebezem, R., & Van Loosdrecht, M. C. (2010). A generalized method for thermodynamic state analysis of
510 environmental systems. *Critical Reviews in Environmental Science and Technology*, 40(1), 1-54.
511 <https://doi.org/10.1080/10643380802000974>

512 Kluyver, T., Ragan-Kelley, B., Pérez, F., Granger, B. E., Bussonnier, M., Frederic, J., ... & Willing, C. (2016). Jupyter
513 Notebooks-a publishing format for reproducible computational workflows. *Elpub*, 87-90. 10.3233/978-1-61499-649-
514 1-87

515 Lehmann, J., Hansel, C.M., Kaiser, C. *et al.* (2020). Persistence of soil organic carbon caused by functional
516 complexity. *Nat. Geosci.* **13**, 529–534 <https://doi.org/10.1038/s41561-020-0612-3198718>

517 Robertson, A. D., Paustian, K., Ogle, S., Wallenstein, M. D., Lugato, E., & Cotrufo, M. F. (2019). Unifying soil
518 organic matter formation and persistence frameworks: the MEMS model. *Biogeosciences*, *16*(6), 1225-1248.
519 <https://doi.org/10.5194/bg-16-1225-2019>

520 Schmidt, M., Torn, M., Abiven, S. et al. (2011) Persistence of soil organic matter as an ecosystem
521 property. *Nature* *478*, 49–56. <https://doi.org/10.1038/nature10386>

522 Stegen, J. C., Garayburu-Caruso, V. A., Danczak, R. E., Goldman, A. E., Renteria, L., Torgeson, J. M., and Wells, J.
523 R.: Hyporheic Zone Respiration is Jointly Constrained by Organic Carbon Concentration and Molecular Richness,
524 EGUSphere [preprint], <https://doi.org/10.5194/egusphere-2022-613>, 2023.

525 Song, H.S., Stegen, J.C., Graham, E.B., Lee, J.Y., Garayburu-Caruso, V.A., Nelson, W.C., Chen, X., Moulton, J.D.
526 and Scheibe, T.D. (2020). Representing organic matter thermodynamics in biogeochemical reactions via substrate-
527 explicit modeling. *Frontiers in microbiology*, *11*, p.531756. <https://doi.org/10.3389/fmicb.2020.531756>

528 Stegen, J.C., Johnson, T., Fredrickson, J.K. *et al.* (2018). Influences of organic carbon speciation on hyporheic
529 corridor biogeochemistry and microbial ecology. *Nat Commun* **9**, 585 <https://doi.org/10.1038/s41467-018-02922-9>

530 Stephanopoulos, George, Aristos A. Aristidou, and Jens Nielsen. (1998) Metabolic engineering: principles and
531 methodologies.

532 Tfaily, M. M., Chu, R. K., Toyoda, J., Tolić, N., Robinson, E. W., Paša-Tolić, L., & Hess, N. J. (2017). Sequential
533 extraction protocol for organic matter from soils and sediments using high resolution mass spectrometry. *Analytica*
534 *chimica acta*, *972*, 54-61.

535 Tolic, N., Liu, Y., Liyu, A., Shen, Y., Tfaily, M. M., Kujawinski, E. B., ... & Hess, N. J. (2017). Formularity: software
536 for automated formula assignment of natural and other organic matter from ultrahigh-resolution mass spectra.
537 *Analytical chemistry*, *89*(23), 12659-12665.

538 Wang, G., W.M. Post and M.A. Mayes (2013). Development of microbial-enzyme-mediated decomposition model
539 parameters through steady-state and dynamics analyses, *Ecological Applications*, *23*(1), 255-272,
540 <https://doi.org/10.1890/12-0681.1>

541 Ward, N.D., Keil, R.G., Medeiros, P.M., Brito, D.C., Cunha, A.C., Dittmar, T., Yager, P.L., Krusche, A.V., Richey,
542 J.E. (2013) Degradation of terrestrially derived macromolecules in the Amazon River. *Nature Geoscience*. *6* (7), 530-
543 533. <https://doi.org/10.1038/ngeo1817>

544 Ward, N.D, Muller, K.A., Chen, X., Zhao, Q., Chu, R., Cheng, Z., Wietsma, T.W., Kukkadapu, R.K. (2023).
545 Interactive Effects of Salinity, Redox State, Soil type, and Colloidal Size Fractionation on Greenhouse Gas Production
546 in Coastal Wetland Soils. ESS Open Archive. <https://doi.org/10.31223/X5FM0N>