

# 1 **ORT: A workflow linking genome-scale metabolic models with reactive transport codes**

2 Rebecca L. Rubinstein<sup>1,\*</sup>, Mikayla A. Borton<sup>2</sup>, Haiyan Zhou<sup>1</sup>, Michael Shaffer<sup>2</sup>, David W. Hoyt<sup>3</sup>,  
3 James Stegen<sup>3</sup>, Christopher S. Henry<sup>4</sup>, Kelly C. Wrighton<sup>2</sup> and Roelof Versteeg<sup>1,\*</sup>

4 <sup>1</sup>Subsurface Insights, LLC., Hanover, NH,

5 <sup>2</sup>Department of Soil and Crop Sciences, Colorado State University, Fort Collins, CO,

6 <sup>3</sup>Pacific Northwest National Laboratory, Richmond, WA,

7 <sup>4</sup>Argonne National Laboratory, Lemont, IL

8 \*To whom correspondence should be addressed.

## 9 **Abstract**

10 **Motivation:** Nutrient and contaminant behavior in the subsurface are governed by multiple  
11 coupled hydrobiogeochemical processes which occur across different temporal and spatial scales.  
12 Accurate description of macroscopic system behavior requires accounting for the effects of  
13 microscopic and especially microbial processes. Microbial processes mediate precipitation and  
14 dissolution and change aqueous geochemistry, all of which impacts macroscopic system behavior.  
15 As ‘omics data describing microbial processes is increasingly affordable and available, novel  
16 methods for using this data quickly and effectively for improved ecosystem models are needed.

17 **Results:** We propose a workflow (‘Omics to Reactive Transport – ORT) for utilizing  
18 metagenomic and environmental data to describe the effect of microbiological processes in  
19 macroscopic reactive transport models. This workflow utilizes and couples two open-source  
20 software packages: KBase (a software platform for systems biology) and PFLOTRAN (a reactive  
21 transport modeling code). We describe the architecture of ORT and demonstrate an  
22 implementation using metagenomic and geochemical data from a river system. Our demonstration

23 uses microbiological drivers of nitrification and denitrification to predict nitrogen cycling patterns  
24 which agree with those provided with generalized stoichiometries. While our example uses data  
25 from a single measurement, our workflow can be applied to spatiotemporal metagenomic datasets  
26 to allow for iterative coupling between KBASE and PFLOTRAN.

27 **Availability and Implementation:** Interactive models available at  
28 <https://pflotranmodeling.paf.subsurfaceinsights.com/pflotran-simple-model/>. Microbiological data  
29 available at NCBI via BioProject ID PRJNA576070. ORT Python code available at  
30 <https://github.com/subsurfaceinsights/ort-kbase-to-pflotran>. KBase narrative available at  
31 <https://narrative.kbase.us/narrative/71260> or static narrative (no login required) at  
32 <https://kbase.us/n/71260/258>

33 **Contact:** [rebecca.rubinstein@subsurfaceinsights.com](mailto:rebecca.rubinstein@subsurfaceinsights.com) or [roelof.versteeg@subsurfaceinsights.com](mailto:roelof.versteeg@subsurfaceinsights.com)

34 **Supplementary information:** Supplementary data are available online.

35

## 36 **1 Introduction**

37 The critical zone (CZ) – the area between the top of the forest canopy and the bottom of  
38 the groundwater table is essential in sustaining life (Guo and Lin, 2016). Being able to understand  
39 and predict critical zone function is essential for both scientific and operational purposes. This  
40 understanding and prediction requires the accurate representation of key hydrobiogeochemical  
41 ecosystem processes which occur and interact in the critical zone. These ecosystem processes  
42 operate at different scales and have different drivers, but at the same time are tightly  
43 interconnected. For instance, while hydrological processes control the movement of water at  
44 macroscopic scales and are driven by groundwater table gradients, precipitation, and  
45 evapotranspiration whereas microbiological processes (Anantharaman *et al.*, 2016; Long *et al.*,  
46 2016) occur at the microbe scale and are driven by microbial populations, soil properties, aqueous

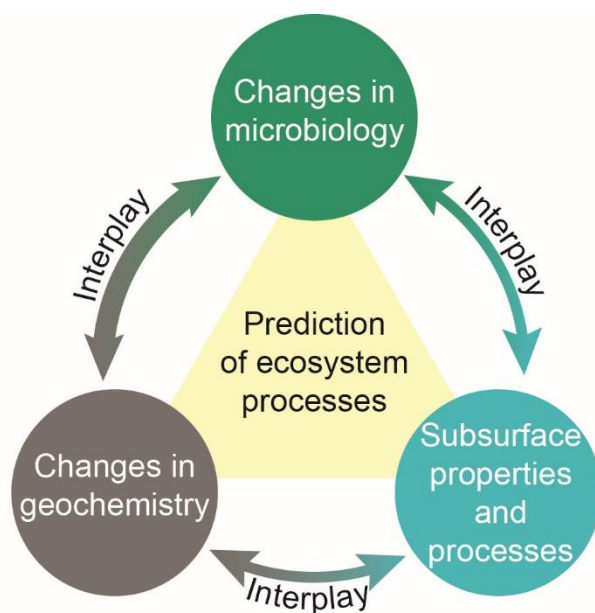
47 geochemistry, and temperature. However, processes at these two scales influence one another in  
48 many ways.

49 One well-established approach to obtaining an understanding of critical zone behavior is  
50 through the use of reactive transport models (RTM) which can simulate coupled chemistry, flow,  
51 and transport in hydrobiogeochemical systems. There are a variety of reactive transport codes (see  
52 (C. I. Steefel *et al.*, 2015) for a review). These models are generally continuum scale models  
53 which represent subsurface properties on grids, with grid volumes on the order of cubic meters. As  
54 the earth is a porous media with grains and pores, such continuum scale models obviously do not  
55 capture pore scale properties and dynamics. One fundamental challenge in numerical modeling is  
56 thus how to link and couple processes and properties which happen at different scales (Battiato *et*  
57 *al.*, 2011; Chu *et al.*, 2012, 2013; Carl I. Steefel *et al.*, 2015). Such linking is especially required  
58 between macroscopic system behavior and microbial processes which change aqueous  
59 geochemistry and mediate precipitation and dissolution.

60 With continued decrease in ‘omics data analysis costs, one promising approach for this  
61 linking is through the incorporation of site-specific microbiological data into RTM to represent  
62 microbe-catalyzed biogeochemical more accurately than using generalized stoichiometries. The  
63 feasibility of using the results of microbiological data analysis to parameterize RTM has been  
64 shown previously. For instance, Scheibe *et al.* demonstrated the linking of genome scale models  
65 with a reactive transport code (in their case, HYDROGEOCHEM) to improve incorporation of  
66 microbiological processes on in situ uranium bioremediation (Scheibe *et al.*, 2009). Specifically,  
67 they used a genome scale model of *Geobacter sulfurreducens* to populate a lookup table spanning  
68 reasonable expected ranges for all combinations of three key system parameters. This was then  
69 used to predict the effects of varying concentrations of three key growth factors (acetate, Fe(III),  
70 and ammonium) on reduction of uranium (VI) at a systems level. More recently, Song *et al.*

71 developed an enzyme-based approach for simulating microbial reaction kinetics which captured  
72 the overall behavior of a consortium rather than rely on individual taxa within the community and  
73 coupled it with reactive transport simulations using PFLOTRAN's Reaction Sandbox (Song and  
74 Liu, 2015; Song *et al.*, 2017; Hammond *et al.*, 2017). This approach is based on a mechanistic  
75 understanding of microbial processes and thus can more accurately predict microbial response to  
76 perturbations. However, this approach substantial experimental data, such as enzyme  
77 concentrations and kinetics data, as well as advanced microbiological knowledge to implement.

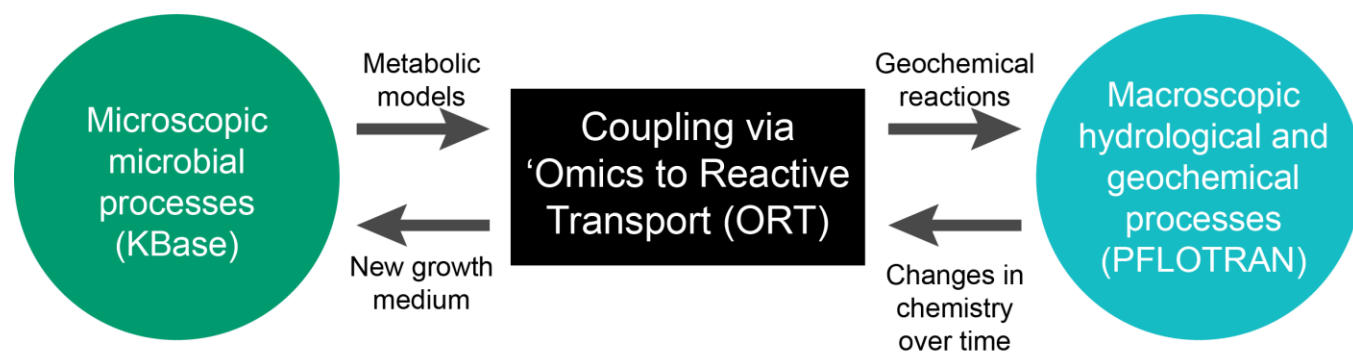
78 These previous efforts have demonstrated the value and feasibility of accurately  
79 representing microbial processes in RTM. This in turn opens up the potential to integrate  
80 microbiological, geochemical, and physical subsurface properties and processes to predict  
81 ecosystem behavior and response (Fig. 1).



82  
83 **Fig. 1** – Accurately capturing the interplay between microbiology, geochemistry, and physical  
84 subsurface properties and processes is critical to understanding and predicting ecosystem  
85 processes.

86           However, each of the approaches described above requires substantial manual effort to  
87   implement for a single site, which makes them challenging to scale. The challenge of scaling these  
88   approaches limits the ability to rapidly develop models obtain the associated understanding for  
89   many sites. An alternative approach (proposed and demonstrated here) is to an approach which  
90   allows for automation.

91           Specifically, in our workflow we use KBase (a cloud-based software platform for systems  
92   biology (Arkin *et al.*, 2018)), to automatically generate draft metabolic models from annotated  
93   metagenome assembled genomes (MAGs) extracted from environmental samples. These  
94   metabolic models can be used (still in KBase) to perform flux balance analysis (FBA) on different  
95   media compositions. These media compositions are informed by metabolomics and other site-  
96   specific chemistry data. The output of the FBA can be used in reactive transport models (such as  
97   the reactive transport model PFLOTRAN (Mills *et al.*, 2009; Hammond and Lichtner, 2010;  
98   Gardner *et al.*, 2015)), which we use in this work. PFLOTRAN is an open source, massively  
99   parallel reactive transport code which supports multi-phase (e.g. aqueous, gaseous), multi-  
100   component (multiple chemical species), and multi-scale (e.g. pore or macroscale) simulation of  
101   contaminant transport in porous media, as well as includes a basic implementation of microbial  
102   reactions modeled by Monod kinetics. One major benefit of PFLOTRAN is that users can  
103   implement custom reactions or kinetics through the Reaction Sandbox (Hammond, 2017). Our  
104   workflow, called ‘Omics to Reactive Transport (ORT) (Fig. 2), thus captures both microbial  
105   metabolisms based on environmental samples (using KBase) and macro-scale hydrologic and  
106   geochemical processes (using PFLOTRAN). In the remainder of this paper we present the concept  
107   and implementation of this workflow.



108

109 **Fig. 2** - Omics to Reactive Transport (ORT) workflow couples microbe-scale and macroscale  
110 processes using the outputs of KBase and PFLOTRAN as inputs for each other.

111

112 In addition to the scientific value of this workflow, we want to highlight three operational  
113 attributes of interest. First, this workflow can be mostly automated, offering the potential of  
114 rapidly generating reactive transport models from microbiological data (detailed in Section 3) with  
115 a minimum of manual labor. Second, the resulting models can easily be shared and made  
116 accessible to other groups. For instance, we have provided two of the models we generated  
117 through a user-friendly web interface which allows end-users to interact with these models. Third,  
118 while in this paper we do not include results for this, our workflow lends itself well to an iterative  
119 approach. Specifically, it is well known that microbial processes will result in changes in  
120 geochemistry, which in turn will influence the microbial processes. In addition to this,  
121 macroscopically driven changes in saturation, temperature, and chemistry (e.g., resulting from  
122 stage-driven surface water/ground water interaction) will also influence microbial processes. The  
123 approach described here can be executed in an iterative manner to capture this two-way coupling  
124 between microscopic and macroscopic processes.

## 125 **2 System & Methods**

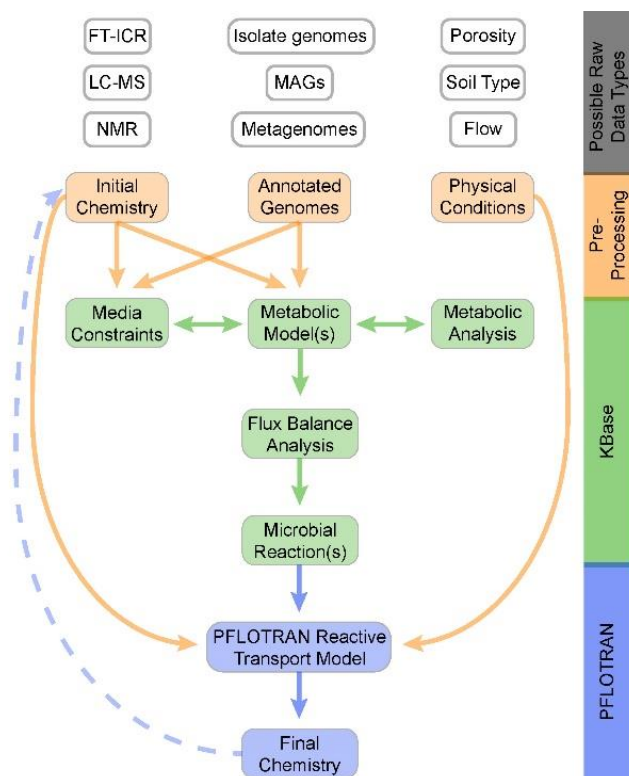
### 126 **2.1 Workflow Concept**

127           The ORT workflow was designed with automation in mind. Specifically, it is designed in a  
128 modular manner with a well-defined start and end points and inputs and outputs, with each  
129 component being fully automatable (Fig. 3). The inputs to this workflow are annotated genomes,  
130 environmental chemistry, and a PFLOTRAN model template which incorporates physical site  
131 data. This template (which would be customized to the specific site) would be something like “0D  
132 batch reactor” or “2D model of unsaturated soil” (where “nD” indicates the number of spatial  
133 dimensions accounted for in the model grid).

134           **In** this workflow we import annotated genomes (Shaffer et al., 2020) and site chemistry  
135 (e.g., available carbon sources, electron acceptors, and micronutrients based on metabolome and  
136 any other chemical analysis at the site, synthesized into a KBase media recipe) into KBase, then  
137 use KBase apps to generate the overall reactions. Next, these reactions, as well as the site  
138 chemistry, physical site data, and model template are used to build the actual PFLOTRAN model.  
139 This model can then be used to simulate macroscopic system behavior.

140           In the iterative implementation, the PFLOTRAN model simulates changes in physical and  
141 chemical conditions in space and time. We can then use the simulated chemistry as a new media  
142 composition to be used by the KBase part of the workflow and repeat the process to generate and  
143 retrieve new resulting overall reactions and substitute them into the PFLOTRAN input file.

144



145

Fig. 3 - Flowchart of the ORT workflow where orange boxes are workflow inputs based on site characterization which are pre-processed before use, green boxes are metabolic modeling steps carried out in KBase, and blue show the resulting RTM. The horizontally-aligned boxes and arrows in the KBase workflow represent robust curation steps (discussed in Section 4), and the dashed arrow indicates the iteration path wherein the PFLORAN-simulated chemistry is used as a new media condition in KBase.

## 146 2.2 Workflow Implementation

147 The ORT workflow consists of Python scripts, KBase narratives, and PFLTORAN models.

148 In our implementation of ORT, KBase apps import genomes and chemical data into KBase and

149 use these as inputs for KBase metabolic modeling apps (process described in detail in Sections 2.4

150 and 2.5). After the completion of the KBase part of the workflow, the KBase API (application

151 programming interface) programmatically exports the KBase-predicted exchange fluxes from

152 KBase. These fluxes are translated by our Python script into an overall reaction string that

153 describes chemical uptake and secretion from each modeled organism, written in PFLORAN-



154 compatible naming conventions. The flux values are used as the stoichiometric coefficients for the  
155 corresponding chemicals in the overall reaction used in PFLOTRAN, with positive fluxes  
156 indicating reactants and negative fluxes indicating products. The summation of exchange fluxes is  
157 not a chemical reaction in the traditional sense, but represents the chemical species removed from  
158 and added to the system as a result of the microbial metabolism. Thus, this “pseudo-reaction”  
159 provides the information needed by PFLOTRAN to simulate the resulting changes in chemical  
160 concentrations.

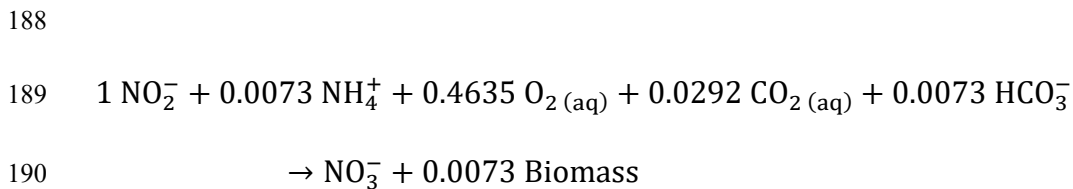
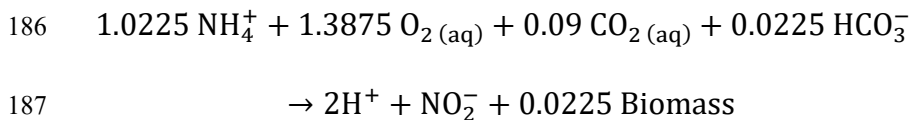
161 The ORT Python script outputs a \*.txt file with the reaction strings and yield terms for use  
162 in the MICROBIAL\_REACTION card in PFLOTRAN as well as a set of \*.dat files which contain  
163 compound names and details which need to be added to the PFLOTRAN geochemical database  
164 (formatted for compatibility with the database). This step can either be done programmatically or  
165 manually by substituting the content of these text files into a PFLOTRAN model input file (known  
166 as an infile). This script bridges the disconnect between KBase and PFLOTRAN illustrated in Fig.  
167 2.

## 168 **2.3 Test Case**

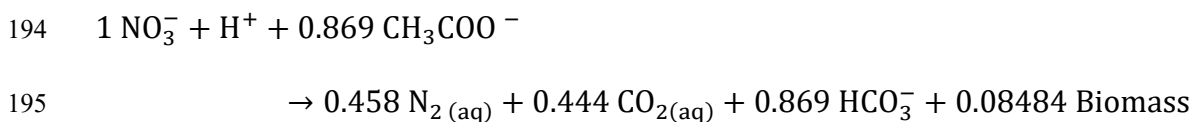
### 169 **2.3.1 System Description**

170 To evaluate the performance of our workflow, we used environmental samples from a  
171 hyporheic zone in the Columbia River. In these zones, biological nitrogen cycling is known to  
172 occur (Triska *et al.*, 1993; Zheng *et al.*, 2016). Biological nitrification and denitrification is a  
173 classic, well-understood, and extensively studied system. We can calculate and compare models  
174 which use traditional (textbook) stoichiometries for nitrification and denitrification to the model  
175 generated from our workflow. In the remainder of this paper, we refer to these two models as the  
176 “literature based model” and “genome derived model”.

177 Nitrification is traditionally split into two sub-processes, ammonium oxidation ( $\text{NH}_4 \rightarrow$   
178  $\text{NO}_2$ ) and nitrite oxidation ( $\text{NO}_2 \rightarrow \text{NO}_3$ ), while denitrification is often represented as a complete  
179 process ( $\text{NO}_3 \rightarrow \text{N}_2$ ), though in reality it is several sequential reactions. Within KBase, we could  
180 implement separate models for each step for which genomes are available, but for comparison to  
181 the traditional model we used a single model for complete denitrification in this test case. The  
182 overall reactions used for the nitrification step were based on experimentally-determined  
183 stoichiometries (Liu and Wang, 2012) determined by fitting data collected from bench-scale  
184 reactors to traditional half-cell reactions (Rittmann and McCarty, 2012), as given by the following  
185 reactions:



191 The complete denitrification process stoichiometry was derived from half-cell reactions  
192 (Rittmann and McCarty, 2012), scaled to one unit nitrate utilization for comparability with the first  
193 two reactions:



196 In both cases, the chemical species represented are limited to classical compositions, which  
197 in some cases may serve as analogs for a range of compounds. These stoichiometries are not  
198 associated with any specific microbes or metabolic pathways, but rather represent the exchange  
199 fluxes observed. While this approach is very effective for process design, it does not offer much  
200 insight into the microbiology of a system, and may obscure finer-scale dynamics – such as less

201 obvious resources that may become limiting or change how the microbes process available  
202 macronutrients, particularly in systems with complex carbon sources.

203 The rates determined through batch kinetics tests (Liu and Wang, 2012) were used for  
204 ammonium oxidation and nitrite oxidation and the denitrification rate was based on rates reported  
205 in the literature (Raboni *et al.*, 2014). The same rates (shown in Table 1) were used for both the  
206 literature-based and genome-based models (described in Section 2.5 - 2.66) in order to directly  
207 compare the effects of the different stoichiometries. In future enhancements, we anticipate that  
208 reaction rates could be used as tunable parameters to fit these models to system-specific  
209 experimental data.

210 Table 1 – Baseline reaction rates used in nitrogen-cycling models

Process	Rate (mol/L·s)	
Ammonium Oxidation	$1.0 \times 10^{-7}$	211 212 213
Nitrite Oxidation	$8.51 \times 10^{-8}$	214
Nitrate Reduction	$2.34 \times 10^{-8}$	215 216

### 217 2.3.2 Leveraging Existing Multiomics Data

218 This study made use of multiomics data from previously published work. Sediment was  
219 collected and DNA extracted as previously described (Graham *et al.*, 2017). To identify the  
220 metabolites available to microorganisms in these river sediments, we performed 1H Nuclear  
221 Magnetic Resonance (NMR) spectroscopy on 17 paired sediment pore water samples which also  
222 had microbial DA extracted, as described previously (Tfaily *et al.*, 2019). Briefly, sediment  
223 samples were mixed with water in a 1:1 ratio and then diluted by 10% (vol/vol) with 5 mM 2,2-  
224 dimethyl-2-silapentane-5-sulfonate-d6 as an internal standard. The 1D 1H NMR spectra of all  
225 samples were processed, assigned, and analyzed using Chenomx NMR Suite 8.3 with  
226 quantification based on spectral intensities relative to the internal standard as described. To obtain  
227 a representative bulk summary of the metabolite environment in these sediments, the

228 concentration of 31 of the NMR identified metabolites was averaged across the 17 sediment  
229 samples, and this data was used as the chemical data input in our ORT workflow (data available in  
230 Supplementary Table S1).

231 Purified genomic DNA was sent to the Joint Genome Institute (JGI, n=33) under  
232 JGI/EMSL proposal 1781 and to the Genomics Shared Resource facility at The Ohio State  
233 University (OSU, n=10), producing 43 metagenomes from 34 sediment samples with an average  
234 sequencing depth of 3.84 (JGI) 25 Gbp (OSU) per sample. JGI and OSU sequencing was  
235 performed as previously described in Graham et al (Graham *et al.*, 2018) and Borton et al (Borton  
236 *et al.*, 2018) respectively. Raw reads were processed, assembled, and binned as outlined in  
237 previous publications (Shaffer *et al.*, 2020) or via the Wrighton Lab GitHub Page  
238 (<https://github.com/TheWrightonLab>). The genomes are available on NCBI via BioProject ID  
239 PRJNA576070.

240 From the sediments, we obtained metagenome assembled genomes (MAGs) from which  
241 we selected four genomes that represented key parts of the nitrogen cycle. For each stage of the  
242 cycle, the most complete genomes capable of filling those roles were selected. To represent  
243 nitrification, we chose the most complete genome representatives of the ammonium oxidizing  
244 archaea classified by GTDB-Tk (version 1.3.0, as of 1-21-21) as a member of the family  
245 Nitrososphaeraceae within the genus TA-21 (previously within the Phylum Thaumarchaeota) and  
246 nitrite oxidizing bacterial member of the Nitrospiraceae for nitrification. Given that the expression  
247 and activity of nitrite reductase encoded in Nitrososphaeraceae (previously Thaumarchaeota) is  
248 poorly understood at this time (Kuypers et al., 2018), we did not incorporate the production of  
249 nitric oxide by Nitrososphaeraceae, and focused only on nitrite outputs from ammonification. To  
250 represent denitrification, we selected two Gammaproteobacterial MAGs, both classified within the  
251 family Steroidobacteraceae. Note that neither of these genomes encoded a gene to produce N<sub>2</sub> gas,

252 but the reaction to convert nitrous oxide to nitrogen gas was added to the metabolic models during  
253 gapfilling (see Section 2.5). We selected only four genomes to maintain the simplicity of this  
254 proof of concept, but the approach could incorporate as many as are needed to capture system  
255 behavior. Each nitrogen-cycling genome was annotated using DRAM (Distilled and Refined  
256 Annotation of Metabolism (Shaffer et al., 2020)) with default parameters. The raw annotations  
257 containing an inventory of all database annotations for every gene from each input genome are  
258 included in the online Supplementary Materials. These genomes and their annotations were  
259 uploaded to KBase (Section 2.5) and were the basis for the KBase-derived model (Section 2.6).

260

## 261 **2.4 Pre-Processing**

262 Prior to executing the workflow, we need to gather and preprocess data and make several  
263 decisions such as selecting a model template. In this section, we describe the data preprocessing  
264 steps in generic terms, as the same steps will be required for any system. To begin our workflow,  
265 user inputs were organized and prepared, which consisted of three broad steps:

266 (1) Qualitative assessment – to balance model complexity and utility, the system definition  
267 phase began with a qualitative description of the system in terms of model type (batch,  
268 chemostat, continuously stirred tank reactor, etc.), important processes (such as nitrification  
269 or sulfur reduction, depending on the system), and parameters of interest (pH, specific  
270 chemical species, etc.) that can guide model development. This step includes evaluating if  
271 there is any “missing” data, which might render the model inaccurate or impossible, and  
272 would need to be estimated in order to produce a viable system (for example, concentrations  
273 of biologically necessary compounds that were not measured). These are identified through a  
274 combination of subject matter knowledge and comparison with KBase default media recipes.  
275 Note that this does not entail delineation of every process and parameter involved in the

276 system, but rather selection of those important to the specific research or application. The  
277 goal of this step is to develop a conceptual model of the system of interest, which may be  
278 augmented and refined as needed to accommodate new data. Because many of these models  
279 will be similar (e.g. 0D batch models), we can build a library of model templates which can  
280 be readily reused.

281 (2) Data Gathering - data describing the site may be drawn from a variety of sources, including  
282 direct sampling at the site and public resources such as weather stations or national  
283 databases. Biological data could come in the form of annotated genomes or metagenomes  
284 collected from the site, or genomes for key microbes as determined using 16S rRNA gene  
285 data or literature review could be drawn from public databases. Chemical data could include  
286 traditional geochemical analysis as well as metabolomics and metaproteomics to provide a  
287 more detailed picture of the chemical profile at the site. Physical data could include  
288 temperature, soil porosity, or other parameters of that nature that would be included in the  
289 PFLOTRAN input file to produce a more site-specific model.

290 (3) Translation to KBase and PFLOTRAN - the data produced by the various analyses above are  
291 not necessarily in formats that may be directly imported to KBase and/or PFLOTRAN.  
292 Therefore, the final step in this phase was to translate these data to forms that can be used by  
293 the relevant tools (KBase or PFLOTRAN). Aside from managing file formats (see the  
294 KBase documentation for details), one major consideration was accounting for any un-  
295 measured chemical species identified in the first step of the preparation phase that needed to  
296 be added to the KBase media composition to make it biologically viable or usable by the  
297 metabolic models generated in KBase. Additions were limited to chemical species or  
298 compounds known (or reasonably expected) to be present and were added in sufficient  
299 concentration that they would not be growth-limiting. The primary check for the presence

300 assumption was that the experimental data indicated that the microbes used were both  
301 present and involved in nitrogen cycling at that site. We did not investigate the assumption  
302 that these compounds were non-limiting, as this is outside the scope of this work.

## 303 **2.5 KBase Metabolic Modeling**

304 Once pre-processing was complete, we can start the ORT workflow. Genomes were uploaded to  
305 KBase as paired FASTA and GFF3 text files using the “Import GFF3/FASTA file as Genome  
306 from Staging Area” app and then annotated with RASTtk using the “Annotate Microbial Genome”  
307 app in KBase. Additional custom annotations from DRAM were uploaded as flat text files using  
308 the beta version of “Import Annotations from Staging” app. If using DRAM annotations,  
309 preprocessing may be carried out using the provided script at  
310 <https://github.com/subsurfaceinsights/ort-kbase-to-pflotran>. Notably, both RASTtk and DRAM are  
311 available as apps in KBase, allowing users to functionally annotate genomes without high memory  
312 computational resources. However, note that the DRAM app in KBase differs from the version  
313 used in this example narrative (Shaffer *et al.*, 2020) as the KBase DRAM app annotates using  
314 KOfam instead of KEGG genes and does not currently include EC reaction identifiers, so end  
315 results may differ from the included narrative. Chemical data was uploaded as flat text files using  
316 the “Import Media file (TSV/Excel) from Staging Area”. The use of pre-processed flat text files as  
317 inputs to the workflow significantly simplifies the process compared to using raw data, especially  
318 for genomes, and these can be generated automatically using scripts such as the one developed for  
319 the DRAM outputs. This first step brought all of our data in the KBase workspace in an integrated  
320 manner.

321 After this step, we used all this data as inputs to the “Build Metabolic Model” app, and the  
322 generated models were used in conjunction with the media objects as inputs to the “Run Flux  
323 Balance Analysis” (FBA) app. Going forward, we refer to this pairing as “growing a model,”

324 meaning we ran the analysis to determine if biomass growth was possible under the given  
325 chemical conditions. The output from the FBA app included the reaction and exchange fluxes for  
326 each model grown on the corresponding media.

327

## 328 **2.6 PFLOTRAN Reactive Transport Modeling**

329 We used our workflow to download the FBA exchange flux values using the KBase API and  
330 translate them from KBase objects with ModelSEED (Henry *et al.*, 2010) compound IDs to flat text  
331 files with reaction strings written using PFLOTRAN naming conventions. We then used either the  
332 KBase-derived reaction strings and biomass yield values or the literature-based stoichiometries  
333 introduced in Section 2.3 to fill in the MICROBIAL\_REACTION card in our OD model template.  
334 All parameters except the reactions and yield terms were held the same for both the literature-based  
335 model and the genome-derived model.

## 336 **3 Model behavior and General behaviors and trends**

337 Both models exhibited sequential ammonium and nitrite oxidation followed by nitrate  
338 reduction, ultimately producing dissolved nitrogen gas (Fig. 4). Despite using the same reaction  
339 rates, inhibition constants, and initial nutrient concentrations, the overall progress of the system is  
340 noticeably different. The genome-derived model exhausts the available ammonium within 1.5 hours  
341 of the simulation start, while the literature based model does not exhaust ammonium until a little  
342 more than 3.5 hours into the simulation. Nitrite concentration peaks earlier and at a lower level for  
343 the genome-derived model (~18  $\mu\text{M}$  at approximately 1 hr) than the literature based model (~51  $\mu\text{M}$   
344 slightly before 3 hrs). Similarly, nitrate peaks at approximately 4  $\mu\text{M}$  after 1.5 hrs for the genome-  
345 derived model but peaks at 40  $\mu\text{M}$  at the 6 hr mark for the literature based model. In the 6 hour  
346 period shown in Fig. 4, the genome derived model has exhausted ammonium, nitrite, and nitrate,  
347 while the literature based model is still processing nitrite and nitrate. This variance is expected since



348 we are comparing generic reactions (with generic substrate utilization and biomass production  
349 reactions) to site-specific reactions based on the most dominant taxa found at our study site.

350 One important difference was that the microbiologically-explicit, genome-based  
351 stoichiometry provided much greater detail on the chemistry, particularly with respect to carbon  
352 catabolism (Fig. 4 and Fig S1). Specifically, the literature-based models relied entirely on either  
353 carbon dioxide (nitrification) or acetate (denitrification), however, because we provided additional  
354 carbon compounds detected from our bulk sediment metabolome, the site models used 15 to 23  
355 unique additional carbon sources, such as betaine, leucine, and choline (see Supplementary Table  
356 1). This greater detail allows us to evaluate more precisely the potential chemical drivers or  
357 limiters of a system which would be entirely overlooked with traditional representations, which  
358 presents the opportunity to probe and improve our conceptual and mechanistic understanding of  
359 these systems and individual metabolisms.

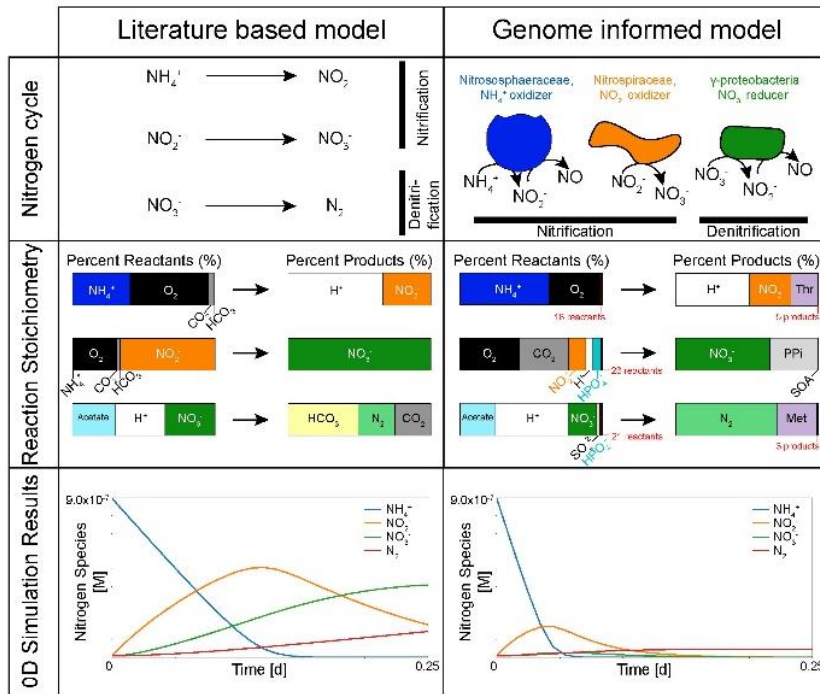


Fig. 4 - Using our Omics to Reactive Transport (ORT) workflow allows us to not only tailor a model to a specific environmental site and system, but also provides much finer insight into the changes in chemistry driven by microbial processes. The top frame shows the steps captured by the literature-based and genome-informed models respectively. The middle frame shows graphical representations of the two sets of reaction stoichiometries. Abbreviations used in the site-specific model frame are Met for Methionine, Thr for Threonine, and SAO for S-Adenosyl-4-methylthio-2-oxobutanoate, which are compounds predicted by KBase as an output which is not part of standard literature representations. The bottom frame shows the results of using each set of reactions in a 0D PFLOTRAN simulation of nitrogen cycling.

360  
 361 Instead of generic bacterial enzymatic reactions, we can determine which site specific bacterial –  
 362 or archaeal – reactions are drivers in the system. Instead of pre-set stoichiometries, our ‘Omics to  
 363 Reactive Transport workflow uses chemistry determined based on metabolomics – using these  
 364 data to describe the initial chemistry rather than generic or simplified chemistry. For example,  
 365 even with the same rate constants, we can see that the genome-informed model utilizes a higher

366 proportion of ammonium in the first step of nitrification, resulting in more rapid depletion of  
367 ammonium in the system and earlier generation of nitrite. As a result, subsequent steps begin  
368 earlier, resulting in an overall accelerated process. At the same time, both versions exhibit the  
369 expected cycling of ammonium to nitrite to nitrate and finally to nitrogen gas. Since PFLOTRAN  
370 relies on user-defined chemistry (as opposed to automatically generating reactions), this allowed  
371 us to incorporate more realistic, mechanism-driven reactions.

372 The genome-based model also allows for greater chemical breadth. The nitrogen cycling  
373 reactions are modulated by a wider range of carbon sources. Additionally, the by-products of this  
374 carbon and nitrogen metabolism also resulted in more complex chemical outputs in some cases,  
375 such as L-Threonine or L-Methionine. These inferred reactions could be further refined by using  
376 gene expression data (e.g., metatranscriptomics or metaproteomics data) to calibrate the models  
377 (by way of reaction rates, saturation constants, etc.) to a particular set of environmental conditions.  
378 Again, this presents an opportunity to test and enhance our understanding of the metabolic  
379 processes involved.

380 Readers can explore and interact with both of these models (without sign-in) through  
381 Subsurface Insights' web-based PFLOTRAN interface at  
382 <https://pflotranmodeling.paf.subsurfaceinsights.com/pflotran-simple-model/>. For the literature-  
383 based model, we have made the input concentrations of ammonium, bicarbonate, and acetate  
384 accessible to web users using sliders. For the Hanford 300 Area-specific version of the model, we  
385 have made accessible the reaction rate for each of the steps modeled. There is no limit to the  
386 number of parameters that may be exposed this way, but for the sake of a user-friendly and un-  
387 cluttered demonstration, we limited our selections to three per model. We selected the parameters  
388 we did both because the effects of varying them are significant and to highlight the power and  
389 flexibility provided by this approach.

## 390 **4 Discussion**

391 We demonstrated an Omics to Reactive Transport (ORT) workflow for creating site  
392 specific reactive transport models that include local chemical and biological content. The ORT  
393 workflow was applied to a well-understood system, and the results agree generally with expected  
394 behavior in a nitrogen cycling system. We interpret the differences in magnitude and timing to be  
395 due to the difference between generic, simplified reactions and metabolism-informed reactions, as  
396 KBase-derived stoichiometries made it possible to capture microbial metabolism in much greater  
397 detail than conventional approaches allow.

398 While the model predictions are borne out by comparison to traditional models, we would  
399 need extensive new data which currently is not available to comprehensively validate our  
400 modeling results. Specifically, we would need high resolution time series data. Such data was not  
401 available in this effort, but is a component of ongoing work, and is in general becoming  
402 increasingly available as technology improves and cost per sample decreases. Given similar data  
403 types, the same workflow could be applied to build and tune a model for other sites.

404 Much of the future work on this workflow will be focused on enhancing and expanding  
405 automation and on making it more robust in several ways. One capability which would be highly  
406 beneficial to our workflow is automated metabolic model curation. In our effort, curation was  
407 carried out manually using two different approaches: metabolism-based and media-based. The  
408 former is labor intensive and requires substantial subject-matter expertise to carry out. The latter is  
409 more straightforward and relies on a more general system understanding, but still requires manual  
410 iteration to obtain reasonable results. Partially or fully automated model curation will eventually  
411 be needed for full automation. This is a topic of active effort by both the KBase core team and  
412 other groups, and we will leverage their efforts. Additional work will be in expanding  
413 PFLOTRAN models to include processes such as temperature mediated biological processes and

414 material recycling. While these are currently not part of the core PFLOTRAN capabilities, these  
415 can be implemented using the PFLOTRAN sandbox.

416 While previous researchers have demonstrated the feasibility of coupling genome-scale  
417 metabolic models with reactive transport simulations, our work is different in some fundamental  
418 ways. First, our workflow, lends itself to automation and rapid model generation from ‘omics data.  
419 As ‘omics data becomes increasingly affordable, the ability to rapidly translate this data into  
420 information on its the implications for macroscopic system behavior will be needed, and our  
421 workflow provides a path towards that. Second, our workflow lends itself to easy incorporation of  
422 more realistic microbial reaction kinetics (e.g., based on temperature or soil conditions). Third, our  
423 workflow lends itself to iteration, which allows us to couple microscopic and macroscopic  
424 processes in either direction. Finally, our workflow provides an easy way to couple two powerful  
425 and complex software packages which typically are used by scientist in different domains, and  
426 allows these scientists a path to generate ‘omics informed reactive transport models.

427

## 428 **Funding**

429 This work has been supported by the SBIR Award DE-SC0019619, Integrated Management and  
430 Analysis Platform for Multi Domain Site Data (program manager Paul Bayer) from the DOE  
431 Biological and Environmental Research program. A portion of the metagenomic sequencing for  
432 this research was performed by the Department of Energy’s Joint Genome Institute (JGI) via  
433 sequencing award no. 1781. Metabolite support was provided by Environmental Molecular  
434 Sciences Laboratory (EMSL) via award no. 50334. Both JGI and EMSL facilities are sponsored  
435 by the Office of Biological and Environmental Research and operated under contract nos. DE-  
436 AC02- 05CH11231 (JGI) and DE-AC05-76RL01830 (EMSL). A portion of this work was  
437 supported by multiple grants within the Wrighton Laboratory: National Sciences Foundation

438 Division of Biological Infrastructure under award no. 1759874, DOE Early Career award no. DE-  
439 SC0018020, and DOE award no. FY21.1068.001. Field sample collection and processing was part  
440 of the Scientific Focus Area (SFA) project at PNNL, sponsored by the U.S. Department of Energy,  
441 Office of Science, Environmental System Science (ESS) Program. This contribution originates  
442 from the ESS Scientific Focus Area (SFA) at the Pacific Northwest National Laboratory (PNNL).

## 443 **Acknowledgements**

444 Tasya Rodzianko, Doug Johnson, and Erek Alper at Subsurface Insights work on the  
445 cyberinfrastructure and web interface that was used in this work. Garret Smith, Pengfei Liu, and  
446 Lindsey Solden provided additional microbiological expertise and processing. Field sample data  
447 was collected by Evan Arntzen, Alex Crump, Brad Fritz, Dave Kennedy, Sarah Fansler, Nate  
448 Phillips, Sadie Montgomery, Kyle Parker, and Rob Macklet at Pacific Northwest National  
449 Laboratory. Processing of fine sediments was also performed by Ray Clayton and Chris Strickland  
450 and cultural support was provided by Doug McFarland and Joy Ferry.

451 *Conflict of Interest:* none declared.

## 452 **References**

- 453 Anantharaman,K. *et al.* (2016) Thousands of microbial genomes shed light on interconnected  
454 biogeochemical processes in an aquifer system. *Nat Commun*, **7**, 13219.  
455 Arkin,A.P. *et al.* (2018) KBase: The United States Department of Energy Systems Biology  
456 Knowledgebase. *Nat Biotechnol*, **36**, 566–569.  
457 Battiato,I. *et al.* (2011) Hybrid models of reactive transport in porous and fractured media.  
458 *Advances in Water Resources*, **34**, 1140–1150.  
459 Borton,M.A. *et al.* (2018) Coupled laboratory and field investigations resolve microbial  
460 interactions that underpin persistence in hydraulically fractured shales. *Proc Natl Acad Sci*  
461 *USA*, **115**, E6585–E6594.  
462 Chu,J. *et al.* (2012) A Multiscale Method Coupling Network and Continuum Models in Porous  
463 Media I: Steady-State Single Phase Flow. *Multiscale Model. Simul.*, **10**, 515–549.  
464 Chu,J. *et al.* (2013) A Multiscale Method Coupling Network and Continuum Models in Porous  
465 Media II—Single- and Two-Phase Flows. In, Melnik,R. and Kotsireas,I.S. (eds), *Advances*  
466 *in Applied Mathematics, Modeling, and Computational Science*, Fields Institute  
467 Communications. Springer US, Boston, MA, pp. 161–185.  
468 Gardner,W.P. *et al.* (2015) High Performance Simulation of Environmental Tracers in  
469 Heterogeneous Domains. *Groundwater*, **53**, 71–80.



- 470 Graham,E.B. *et al.* (2018) Multi 'omics comparison reveals metabolome biochemistry, not  
471 microbiome composition or gene expression, corresponds to elevated biogeochemical  
472 function in the hyporheic zone. *Science of The Total Environment*, **642**, 742–753.
- 473 Guo,L. and Lin,H. (2016) Critical Zone Research and Observatories: Current Status and Future  
474 Perspectives. *Vadose Zone Journal*, **15**.
- 475 Hammond,G.E. *et al.* (2017) Application of a hybrid multiscale approach to simulate hydrologic  
476 and biogeochemical processes in the river-groundwater interaction zone. Sandia National  
477 Lab. (SNL-NM), Albuquerque, NM (United States).
- 478 Hammond,G.E. (2017) PFLOTRAN Reaction Sandbox: A Flexible Extensible Framework for  
479 Vetting Biogeochemical Reactions within an Open Source Subsurface Simulator.
- 480 Hammond,G.E. and Lichtner,P.C. (2010) Field-scale model for the natural attenuation of uranium  
481 at the Hanford 300 Area using high-performance computing: MODEL FOR NATURAL  
482 ATTENUATION OF URANIUM. *Water Resour. Res.*, **46**.
- 483 Henry,C.S. *et al.* (2010) High-throughput generation, optimization and analysis of genome-scale  
484 metabolic models. *Nature Biotechnology*, **28**, 977–982.
- 485 Kuypers,M.M.M. *et al.* (2018) The microbial nitrogen-cycling network. *Nat Rev Microbiol*, **16**,  
486 263–276.
- 487 Liu,G. and Wang,J. (2012) Probing the stoichiometry of the nitrification process using the  
488 respirometric approach. *Water Research*, **46**, 5954–5962.
- 489 Long,P.E. *et al.* (2016) Microbial Metagenomics Reveals Climate-Relevant Subsurface  
490 Biogeochemical Processes. *Trends in Microbiology*, **24**, 600–610.
- 491 Mills,R.T. *et al.* (2009) Modeling subsurface reactive flows using leadership-class computing. *J.*  
492 *Phys.: Conf. Ser.*, **180**, 012062.
- 493 Raboni,M. *et al.* (2014) Calculating specific denitrification rates in pre-denitrification by assessing  
494 the influence of dissolved oxygen, sludge loading and mixed-liquor recycle. *Environmental*  
495 *Technology*, **35**, 2582–2588.
- 496 Rittmann,B.E. and McCarty,P.L. (2012) Environmental biotechnology: principles and applications  
497 Tata McGraw-Hill Education.
- 498 Scheibe,T.D. *et al.* (2009) Coupling a genome-scale metabolic model with a reactive transport  
499 model to describe in situ uranium bioremediation. *Microbial Biotechnology*, **2**, 274–286.
- 500 Shaffer,M. *et al.* (2020) DRAM for distilling microbial metabolism to automate the curation of  
501 microbiome function. *Nucleic Acids Res*, **48**, 8883–8900.
- 502 Song,H.-S. *et al.* (2017) Regulation-Structured Dynamic Metabolic Model Provides a Potential  
503 Mechanism for Delayed Enzyme Response in Denitrification Process. *Frontiers in*  
504 *Microbiology*, **8**.
- 505 Song,H.-S. and Liu,C. (2015) Dynamic Metabolic Modeling of Denitrifying Bacterial Growth:  
506 The Cybernetic Approach. *Industrial & Engineering Chemistry Research*, **54**, 10221–  
507 10227.
- 508 Steefel,Carl I. *et al.* (2015) Micro-Continuum Approaches for Modeling Pore-Scale Geochemical  
509 Processes. *Reviews in Mineralogy and Geochemistry*, **80**, 217–246.
- 510 Steefel,C. I. *et al.* (2015) Reactive transport codes for subsurface environmental simulation.  
511 *Computational Geosciences*, **19**, 445–478.
- 512 Tfaily,M.M. *et al.* (2019) Single-throughput Complementary High-resolution Analytical  
513 Techniques for Characterizing Complex Natural Organic Matter Mixtures. *JoVE*, 59035.
- 514 Triska,F.J. *et al.* (1993) The role of water exchange between a stream channel and its hyporheic  
515 zone in nitrogen cycling at the terrestrial—aquatic interface. In, Hillbricht-Ilkowska,A. and  
516 Pieczyńska,E. (eds), *Nutrient Dynamics and Retention in Land/Water Ecotones of*  
517 *Lowland, Temperate Lakes and Rivers*, Developments in Hydrobiology. Springer  
518 Netherlands, Dordrecht, pp. 167–184.
- 519 Villa,J.A. *et al.* (2020) Methane and nitrous oxide porewater concentrations and surface fluxes of a  
520 regulated river. *Science of The Total Environment*, **715**, 136920.
- 521 Zheng,L. *et al.* (2016) Temperature effects on nitrogen cycling and nitrate removal-production  
522 efficiency in bed form-induced hyporheic zones. *Journal of Geophysical Research:*  
523 *Biogeosciences*, **121**, 1086–1103.