

**An Artificial Intelligence Based Self-Adaptive Dynamic Process Control System for  
Enhancing In-Situ Bioremediation of Benzene-Contaminated Groundwater**

Li He<sup>1\*†</sup>, Xu Duan<sup>1†</sup>, Chenyang Li<sup>1</sup>, Xu Yang<sup>1</sup>, Mengxi He<sup>1</sup>

<sup>1</sup>State Key Laboratory of Hydraulic Engineering Simulation and Safety, Tianjin University,  
Yaguan Road 135, Haihe Education Park, Jinnan District, Tianjin, 300350, China

\*Corresponding author: Li He (helix111@tju.edu.cn)

## Contents

<b>Text S1. Set-up of the pilot-scale reactor</b> .....	3
<b>Text S2. Pilot-scale experimentation analysis</b> .....	5
<b>Text S3. Contaminant fate and transport modeling in the groundwater</b> .....	7
<b>Text S4. Biodegradation modeling of contaminants in the groundwater</b> .....	9
<b>Text S5. Procedures for solving the coupled flow and transport problem</b> .....	14
<b>Text S6. Stepwise cluster analysis (SCA)</b> .....	15
<b>Text S7. Filtering Process Model</b> .....	18
<b>Text S8. Nonlinear Optimization model of the FCI optimizer</b> .....	21
<b>Text S9. Nonlinear Optimization model of the DPC system</b> .....	22
<b>Text S10. Nonlinear Optimization model of the SADPC system</b> .....	23
<b>Text S11. Genetic algorithms (GA)</b> .....	24
<b>Text S12. MPC control module procedure</b> .....	25
<b>Text S13. General procedure for developing a process control system for enhanced in situ biodegradation</b> .....	26
<b>Text S14. Collinearity test of the independent variables</b> .....	27
<b>Figure S1(a). Plan view of the pilot scale system</b> .....	28
<b>Figure S1(b). Front view of the pilot scale system</b> .....	28
<b>Figure S1(c). Bottom view of the pilot scale system</b> .....	29
<b>Figure S1(d). End elevation of the pilot scale system</b> .....	29
<b>Figure S1(e). Well locations (plan view)</b> .....	30
<b>Figure S1(f). Well locations and soil types (section view)</b> .....	31
<b>Figure S2. Framework of the FCI Optimizer</b> .....	31
<b>Figure S3. Framework of the SI Emulator</b> .....	32
<b>Figure S4. Locations of the hypothetical wells</b> .....	32
<b>Figure S5. Benzene concentrations of the DPC system from Day 2 to Day 22, where Figs. (a) to (i) represents the concentrations at HW-40, HW-42, HW-48, HW-52, HW-56, HW-62, HW-95, HW-102, and HW-106</b> .....	33
<b>Figure S6. The concentration distribution of Benzene on Day 57 of the experiment</b> .....	33
<b>Figure S7(a). The concentration distribution of Benzene on Day 40</b> .....	34
<b>Figure S7(b). The concentration distribution of Benzene on Day 57</b> .....	34
<b>Figure S8. Verification results for well 5 and well 6</b> .....	35
<b>Table S1. Input parameters for contaminant transport simulation</b> .....	36
<b>Table S2. Initial Geochemical and Microbial Properties of the Soil</b> .....	37
<b>Table S3. Observed benzene concentrations (mg/L)</b> .....	38
<b>Table S4. Error analysis for the biodegradation simulation results</b> .....	40
<b>Table S5. Fifty levels of contamination situation (mg/L)</b> .....	41
<b>Table S6. Fifty scenarios of operating conditions</b> .....	42
<b>Table S7. Input and output variables for SI emulator and FCI simulator</b> .....	43
<b>Table S8. The result of the collinearity test</b> .....	44

### Text S1. Set-up of the pilot-scale reactor

The modeling domain is defined as three-dimensional (3-D), with the contaminated zone around the groundwater table considered as the major pollution source. The area of the simulation domain is  $3.6 \times 1.2 \text{ m}^2$ . Vertically, the simulation domain is discretized into four grid blocks corresponding to four simulation layers; each layer is located in the middle of the grid block that facilitates the application of a block-centered finite difference scheme. In the horizontal plain, each layer is discretized into  $24 \times 8$  grids. Each grid has dimensions of 0.15, 0.15, and 0.30 m in x, y, and z directions, respectively. The total number of grids in this 3-D computational system is 768 ( $24 \times 8 \times 4$ ). Layers 3 and 4 are located in the saturated zone, while layers 1 and 2 are situated in the unsaturated zone. The detailed views of the pilot-scale reactor can be seen in Figure S1(a)-(d) of the supporting information.

Three soil distributions in the four layers are shown in Figure S1(f). The monitoring wells were used to obtain the subsurface hydrogeology “representative” view (well locations are shown in Figure S1(e)). The LNAPLs (Light Non-Aqueous Phase Liquids) initially occupied a contaminated area in layers 3 and 4 around the groundwater table. The zero-flow boundary conditions were enforced at the top and bottom of the modeling domain, as well as at the sides parallel to the x-axis. Constant hydraulic heads were employed at the left and right boundaries, allowing continuous water flow in the aquifer. Benzene concentrations in the system can be forecasted through the developed simulator. Table S1 presents the input parameters for the simulation model.

Water-table level gauges were installed in the first and fourth sections to monitor water depths inside the reactor. Observation windows were built on the front side of each section while another on the top. The side windows were used to observe the subsurface conditions, and the top ones to observe the soil surface. The four sections were connected to each other with flanges, each of which had 44 bolts. Gaskets made of anti-organic solvent and anti-high temperature rubber and silicone pastern were placed between the flanges to prevent the leakages.

For soils loading, the pre-selected clay was from the construction site of the Saskatchewan Indian Federated College, Regina, Saskatchewan, at depths of 2 to 6 m from the ground surface. The clayey till and fine sand were provided by the Waxy’s Bobcat and Landscaping Ltd. The initial properties of soils can be seen in Table S2. The value of soil organic carbon (SOC) in the soils is measured at around 1.14%. Some activities around the concentration site accelerate the decomposition of the organic carbon, and there is little vegetation and insufficient organic carbon input, resulting in the particularly low concentration. Therefore, the transport of organic carbon is not considered in this model, which is assumed to migrate with the movement of the organic phase. We thus ignored the effect of SOC in this study, also considering the high concentration of benzene contamination, as treated by many existing studies (Jimenez et al., 2006; Wolicka et al., 2009; Xin et al., 2013; Umar et al., 2021; Yang et al., 2019). For instance, Liang et al (2013) mentioned that benzene can be used as the sole source of organic carbons when discussing the chlorobenzene and benzene degradation in the groundwater. Besides we did not directly consider the effect of the grain size, instead of using porosity and intrinsic permeability when describing the transport of the target contaminant in the mathematical model. We did it by hypothesizing the grain size can affect soil porosity and intrinsic permeability and have a further impact on the fate and transport of benzene. This method is generally used in many groundwater models such as MODFLOW and NAPL3D (McDonald & Harbaugh, 1988; Zhang et al., 2008;

Niswonger & Prudic, 2013; Hu et al., 2021). Detailed information and guides can be seen in USGS website (<https://www.usgs.gov/mission-areas/water-resources/science/modflow-and-related-programs>) and EPA website (<https://www.epa.gov/water-research/non-aqueous-phase-liquid-napl-simulator>).

The inner surface of the reactor wall was labeled and divided into grids where different types of soils were loaded. Clay and clay-till were sieved by 1/4 inch sieving meshwork before loading. Sand was firstly loaded followed by clay-till and then clay. Upon the completion of every 100 mm depth in every grid, tap water was spread on the surface, and the soil in the grid was vibrated by concrete vibrator and pressed by impinging a hammer on a wood board that directly contacted the soil surface to ensure homogeneity and non-fracture structure. Upon the completion of each layer, more water was spread and the layer was left overnight to settle. After the soil loading process was completed, the system was then left still for three months with a water flow of 10 L/day through the loaded soils. No noticeable further settlement or consolidation was observed afterwards.

A thermostatic room, in which the pilot-scale system and the accessorial equipment were assembled, was built to realize various temperatures by an air conditioner. Water and drainage containers were each connected to the upstream inlet and downstream outlet, respectively. Water level gauges were used to show the depth of water table in the reactor. Tap water in a water container was pumped into the reactor through six water inlets on the inlet-end board as upstream groundwater inflow through a peristaltic pump. Before the start of the experiments, water in the container was kept still overnight to reach the room temperature. A 7-day buffering time preceded all experiments in the pilot model so that the temperature at every location in the reactor could reach equilibrium. The upstream water was kept flowing through system to acquire the desired soil temperature.

The monitoring wells are for facilitating access to the groundwater so that a “representative” view of the subsurface hydrogeology can be obtained, either through the collection of water samples or the measurement of physical and hydraulic parameters. In this study, a few monitoring wells were also used for pumping and injecting purposes during the remediation processes. Locations of the wells are presented in Figure S1(e). There are 25 wells allocated in four sections of the pilot system. Soil in the system was stratified into four layers, with the third and fourth layers being saturated with water. Each layer is 30 cm deep. Among the wells, 13 of them (with PVC pipes) were installed to reach the third soil layer; the other 12 wells could reach the fourth layer (Figure S1(f)). Small holes were uniformly made around the bottom sections of the pipes. Screens were used to wrap the pipes to prevent from soil clogging. Soil particles were prevented from moving into the wells while the groundwater could infiltrate into them. The wells were sealed by rubber caps at the tops. For each well, a hose was installed that passed through the caps and reached its bottom. The outside of the hose was clamped by a clip so that air and groundwater in the well were isolated from the atmosphere.

## Text S2. Pilot-scale experimentation analysis

Benzene concentrations monitored in both natural-attenuation and bioremediation phases are listed in Table S3 and Figure S6(a-b) of the Supporting Information. The highest benzene concentrations were encountered in well 6 during the entire period of the natural-attenuation phase. This was due to the fact that well 6 was close to the leakage source. High concentrations were also observed in wells 3 and 10, which were placed in the third layer. The contaminant can easily reach these wells along with groundwater flow since the leakage occurred at the top of the third layer. Moreover, much of the contaminant transported within the third layer and did not migrate to the fourth one since gasoline is lighter than water. In comparison, the highest concentrations in the fourth layer were observed in well 5, which was installed in the sand zone near the leakage source. Due to the low porosity and permeability of silty and clayey soils, benzene was not observed in the down gradient domain of the pilot system until day 32 (the contaminant reached well 16 on day 32).

On day 40, enhanced in-situ biodegradation action was undertaken. It is shown from Table S3 of the Supporting Information that the benzene concentrations vary greatly due to flow-condition changes resulting from the pumping and injecting actions. The location of the peak concentration moved towards the downstream. The benzene concentrations in the groundwater also decreased greatly compared with those in the earlier periods. The peak benzene concentration decreased from 7.34 mg/L at the beginning of the remediation program to 0.633 mg/L on day 17 after the remediation action started. It is indicated that the enhanced in-situ bioremediation had efficiency in removing benzene from the groundwater.

The experimental results indicated that the developed pilot-scale reactor can effectively facilitate the simulation of both natural attenuation and enhanced remediation processes. The experimental results can be used for validating, calibrating and verifying the developed numerical model under different site conditions.

Calibration and verification of the developed biodegradation model were undertaken using data obtained from the pilot experiments. The results of the error analysis are provided in Table S4 of the Supporting Information. The absolute errors between the simulated and observed concentrations of 12 wells range from 0.00 to 0.40 mg/L with a mean of 0.21 mg/L. The root-mean-square error is 0.27 mg/L, and the correlation coefficient is 0.93. According to the error analysis, the biodegradation simulation result has been proved to be within a relatively reasonable range. Figure S7 shows the verification results for Day 57, which is the end of phase I bioremediation. The concentration distributions of benzene are generally the same based on observed data and simulated data. The highest concentration levels are obtained at the bottom right corner of the experiment region, and the coordinates of these points are around (2.1, 0.6) and (3, 0.45). The verification results for the temporal variations of benzene concentrations in well 5 and well 6 are shown in Figure S8, which indicates that this

model can simulate the actual degradation process of benzene properly. After the calibration and verification, the simulation model can be used for investigating the effects of different bioremediation strategies on benzene concentrations.

### Text S3. Contaminant fate and transport modeling in the groundwater

A critical step in understanding the impact of a subsurface release of NAPL is a modeling analysis of the NAPL flow and transport and fate of its crucial constituents. A complete description of multiphase flow and transport in subsurface must include flow of the fluid phases (water, gas, NAPL, etc.), mass transfer of species between these phases, and transport of species in each phase. A three-dimensional multiphase multicomponent SEAR model is recognized as an effective tool in investigating complex physical processes involved in NAPLs flow and transport. Several models have been developed to simulate the flow of multiple fluid phases in subsurface during recent years. All these models included simplifying assumptions with respect to phase presence and dimensionality, NAPL contaminant mass balance and water and air in subsurface. Important assumptions used in the development of mass conservation equations are: 1) the solid phase is immobile; 2) soil and fluids are slightly compressible; 3) dispersion is of Fickian form; 4) components mix ideally; 5) Darcy's law applies in the calculation of phase velocities. The basic mass conservation equation for components in subsurface can be written as (Delshad et al., 1996):

$$\frac{\partial}{\partial t}(\phi \tilde{C}_k \rho_k) + \vec{\nabla} \cdot \left[ \sum_{l=1}^{n_p} \rho_k (C_{kl} \vec{u}_l - \phi S_l \vec{\bar{D}}_{kl} \cdot \vec{\nabla} C_{kl}) \right] = R_k \quad (S1)$$

where  $k$  is the component index;  $l$  is the phase index;  $\phi$  is the soil porosity;  $\tilde{C}_k$  is the overall concentration of component  $k$  (volume fraction);  $\rho_k$  is the density of component  $k$  [ML<sup>-3</sup>];  $n_p$  is the number of phases;  $C_{kl}$  is the concentration of component  $k$  in phase  $l$  (volume fraction);  $\vec{u}_l$  is the Darcy velocity of phase  $l$  [LT<sup>-1</sup>];  $S_l$  is the saturation of phase  $l$ ;  $R_k$  is the total source/sink term for component  $k$  (volume of the component  $k$  per unit volume of porous media per unit time);  $\vec{\bar{D}}_{kl}$  is the dispersion tensor. The overall concentration ( $\tilde{C}_k$ ) denotes the volume of the component  $k$  summed over all phases.

The dispersion tensor ( $\vec{\bar{D}}_{kl}$ ) can be expressed as follows (Bear, 1979):

$$\vec{\bar{D}}_{kl} = \frac{D_{m,kl}}{\tau} \delta_{ij} + \frac{\alpha_{Tl}}{\phi S_l} |\vec{u}_l| \delta_{ij} + \frac{(\alpha_{Ll} - \alpha_{Tl}) u_{li} u_{lj}}{\phi S_l |\vec{u}_l|} \quad (S2)$$

where  $\tau$  is tortuosity (defined as a value greater than 1);  $D_{m,kl}$  is molecular diffusion coefficient of component  $k$  in phase  $l$  [L<sup>2</sup>T<sup>-1</sup>];  $\delta_{ij}$  is Kronecker delta function;  $\alpha_{Ll}$  and  $\alpha_{Tl}$  are longitudinal and transverse dispersivities of phase  $l$ , respectively [L];  $u_{li}$  and  $u_{lj}$  are Darcy velocities of phase  $l$  in directions  $i$  and  $j$ , respectively [LT<sup>-1</sup>];  $|\vec{u}_l|$  is magnitude of the vector flux for phase  $l$  [LT<sup>-1</sup>]. The phase flux can be calculated from the multiphase form of the Darcy's law (Faust et al., 1989):

$$\vec{u}_l = -\frac{k_{rl} \vec{\bar{K}}}{\mu_l} \cdot (\vec{\nabla} P_l - \rho_l g \vec{\nabla} z) \quad (S3)$$

where  $k_{rl}$  is relative permeability of porous medium to phase  $l$ ;  $\vec{\bar{K}}$  is intrinsic permeability tensor [L<sup>2</sup>];  $\mu_l$  is viscosity of phase  $l$  [ML<sup>-2</sup>T<sup>-1</sup>];  $\rho_l$  is density of phase  $l$  [ML<sup>-3</sup>];  $g$  is acceleration of gravity [LT<sup>-2</sup>];  $z$  is vertical distance which is defined as positive downward [L];  $P_l$  is pressure of phase  $l$  [ML<sup>-1</sup>T<sup>-2</sup>].

Grumberg and Nissan's correlation is used to calculate the NAPL viscosity as a function of organic species concentration:

$$\ln \mu_{2l} = \sum_{k=1}^{n_o} x_{kl}^o \ln \mu_k^o \quad (\text{S4})$$

where  $\mu_{2l}$  is the organic mixture viscosity;  $x_{kl}^o$  is molar fraction of each organic component in phase  $l$  (water, NAPL, etc.);  $\mu_k^o$  is the viscosity of single organic component; and  $n_o$  is the number of organic components in NAPL. For a NAPL mixture, the overall organic hydrostatic pressure gradient is obtained by assuming ideal mixing:

$$C_{2l} \gamma_{2l} = \sum_{k=1}^{n_o} C_{kl}^o \gamma_{kl}^o \quad (\text{S5})$$

where  $C_{2l}$  is concentration of NAPL mixture in phase  $l$  (water, NAPL, etc.);  $\gamma_{2l}$  is density of NAPL mixture in phase  $l$ ;  $C_{kl}^o$  is concentration of each organic component in phase  $l$ ;  $\gamma_{kl}^o$  is density of single organic component in phase  $l$ ; and  $n_o$  is the number of organic components in NAPL.

The aquifer boundaries are modeled as either constant potential surfaces or closed surfaces. The model can be solved numerically through the block-centered finite difference method. The solution method for the contaminant-transport model is the implicit pressure-explicit saturation method. The only unknown in the pressure equation is the pressure of water phase.



#### Text S4. Biodegradation modeling of contaminants in the groundwater

Generally, the biodegradation model involves simulation of substrate competition, nutrient limitation, product toxic inhibition, and aerobic cometabolism. The basic structure of the biodegradation model for a system with single substrate, single electron acceptor and single biological species can be characterized as follows (de Blanc, 1998):

$$\frac{dC^{APS}}{dt} = -\frac{A^{SM}k^{SMT}\overline{C^{AAB}}}{m^{CSM}}(C^{APS} - \overline{C^{ABS}}) \quad (S6)$$

$$\begin{aligned} \frac{d\overline{C^{ABS}}}{dt} = & \frac{A^{SM}k^{SMT}}{V^{SM}}(C^{APS} - \overline{C^{ABS}}) - \frac{\mu_{max}C^{AUB}}{k^Y}\left(\frac{C^{APS}}{k^{SHS} + C^{APS}}\right)\left(\frac{C^{APE}}{k^{EHS} + C^{APE}}\right) - k^{FRR}C^{APS} \\ & - k^{FRR}\overline{C^{ABS}} \end{aligned} \quad (S7)$$

$$\frac{dC^{APE}}{dt} = -\frac{A^{SM}k^{EMT}\overline{C^{AAB}}}{m^{CSM}}(C^{APE} - \overline{C^{ABE}}) \quad (S8)$$

$$\begin{aligned} \frac{d\overline{C^{ABE}}}{dt} = & \frac{A^{SM}k^{EMT}}{V^{SM}}(C^{APE} - \overline{C^{ABE}}) - \frac{\mu_{max}C^{AUB}m^{EAC}}{k^Y}\left(\frac{C^{APS}}{k^{SHS} + C^{APS}}\right)\left(\frac{C^{APE}}{k^{EHS} + C^{APE}}\right) \\ & - \frac{\mu_{max}\rho_X m^{EAC}}{k^Y}\left(\frac{\overline{C^{ABS}}}{k^{SHS} + \overline{C^{ABS}}}\right)\left(\frac{\overline{C^{ABE}}}{k^{EHS} + \overline{C^{ABE}}}\right) \end{aligned} \quad (S9)$$

$$\frac{dC^{AUB}}{dt} = \mu_{max}C^{AUB}\left(\frac{C^{APS}}{k^{SHS} + C^{APS}}\right)\left(\frac{C^{APE}}{k^{EHS} + C^{APE}}\right) - k^{ED}C^{AUB} \quad (S10)$$

$$\frac{d\overline{C^{AAB}}}{dt} = \mu_{max}\overline{C^{AAB}}\left(\frac{\overline{C^{ABS}}}{k^{SHS} + \overline{C^{ABS}}}\right)\left(\frac{\overline{C^{ABE}}}{k^{EHS} + \overline{C^{ABE}}}\right) - k^{ED}\overline{C^{AAB}} \quad (S11)$$

where  $C^{APS}$  is the aqueous phase substrate concentration (substrate mass per unit volume of aqueous phase);  $\overline{C^{ABS}}$  is the substrate concentration in attached biomass (mass of substrate per unit volume of biomass);  $C^{APE}$  is the aqueous phase electron acceptor concentration (mass of electron acceptor per unit volume of aqueous phase);  $\overline{C^{ABE}}$  is the electron acceptor concentration in attached biomass (mass of electron acceptor per unit volume of biomass);  $C^{AUB}$  is the aqueous phase concentration of unattached biomass (mass of unattached cells per unit volume of aqueous phase);  $\overline{C^{AAB}}$  is the attached biomass concentration (mass of attached cells per volume of aqueous phase);  $A^{SM}$  is the surface area of a single microcolony ( $L^2$ );  $k^{EMT}$  is the electron acceptor mass transfer coefficient ( $LT^{-1}$ );  $k^{SMT}$  is the substrate mass transfer coefficient ( $LT^{-1}$ );  $\mu_{max}$  is the maximum specific growth rate ( $T^{-1}$ );  $m^{CSM}$  is the mass of cells in a

single microcolony,  $m^{\text{CSM}} = \rho V^{\text{SM}}$  (M);  $m^{\text{EAC}}$  is the mass of electron acceptor consumed per mass of substrate biodegraded;  $\rho$  is the biomass density (mass of cells per volume of biomass);  $V^{\text{SM}}$  is the volume of a single microcolony ( $\text{L}^3$ );  $k^{\text{Y}}$  is the yield coefficient (mass of cells per volume of biomass);  $k^{\text{SHS}}$  is the substrate half-saturation coefficient ( $\text{ML}^{-3}$ );  $k^{\text{EHS}}$  is the electron acceptor half-saturation coefficient ( $\text{ML}^{-3}$ );  $k^{\text{FRR}}$  is the first-order reaction rate coefficient (for abiotic decay reactions,  $\text{T}^{-1}$ );  $k^{\text{ED}}$  is the endogenous decay coefficient ( $\text{T}^{-1}$ ); and  $t$  is the time (T).

Reduction of contaminants in the aqueous phase in Equation (S6) results from three mechanisms. The first term accounts for diffusion of contaminants from liquid phase across a stagnant liquid film into attached biomass. The second one indicates the reduction of contaminants by unattached microorganisms in the bulk liquid. The reduction rate is affected by concentrations of contaminant and electron acceptor through the Monod kinetic. Substrate competition, nutrient limitations, inhibition, and reducing power limitations can also be incorporated within the second term as described in the following sections. The third term accounts for abiotic loss of contaminants through first-order reactions. One equation of the same form as Equation (S6) will be used for each substrate.

Equation (S7) describes the loss of substrate within attached biomass. It describes processes of substrate diffusion into attached biomass, biodegradation within the biomass, and abiotic decay. Substrate competition, nutrient limitations, inhibition, and reducing power limitations can also be incorporated into this term for biodegradation of the substrate. Equations (S8) and (S9) describe the loss of the electron acceptor, which are of the same form as Equations (S6) and (S7). Equations (S8) and (S9) simulate the growth and decay of unattached and attached biomass, respectively.

The attached biomass concentration ( $\overline{C^{\text{AAB}}}$ ) is dependent upon the biomass density, microcolony volume and microcolony mass (de Blanc, 1998):

$$\overline{C^{\text{AAB}}} = \frac{N^{\text{CPS}} D^{\text{PM}} m^{\text{CSM}}}{N^{\text{CPM}} \phi} \quad (\text{S12})$$

where  $N^{\text{CPS}}$  is the number of cells per mass of solid;  $D^{\text{PM}}$  is the bulk density of the porous medium;  $N^{\text{CPM}}$  is the number of cells per microcolony (a constant); and  $\phi$  is the porosity.

Since the biomass density, number of cells, mass of one microcolony, and medium porosity are assumed to be constant,  $\overline{C^{\text{AAB}}}$  is proportional to  $N^{\text{CPS}}$  or, alternately, to  $(N^{\text{CPS}} / N^{\text{CPM}})$ , (the number of microcolonies). Moreover, the area available for transport of species from the aqueous phase to the biomass is directly proportional to  $\overline{C^{\text{AAB}}}$ , because the surface area per microcolony is assumed constant.

### Multiplicative Monod kinetics

For multiplicative Monod kinetics, it is assumed that other limiting nutrients are also limiting microbial growth besides substrates and electron acceptors. When other

chemical species or nutrients such as nitrogen or phosphorous are limiting factors, the substrate utilization term can be modified correspondingly in order to account for these additional limitations (Rittmann et al., 1991):

$$\mu'_{max} = \mu_{max} \cdot \frac{C^{LN}}{k^{LNH} + C^{LN}} \quad (S13)$$

where  $C^{LN}$  is concentration of a limiting nutrient ( $ML^{-3}$ ); and  $k^{LNH}$  is limiting nutrient half-saturation coefficient concentration ( $ML^{-3}$ ).

### Biomass growth

The basic biomass growth expression of equations (S10) and (S11) contains an additional term to limit the volume of the biomass. With this limitation, the general form of the biomass growth expression is (de Blanc, 1998):

$$\frac{dC^{AUB}}{dt} = \mu_{max} C^{AUB} \left( \frac{C^{APS}}{k^{SHS} + C^{APS}} \right) \left( \frac{A^{SM}}{K_A + A^{SM}} \right) \left( 1 - \frac{C^{AUB}}{0.9\rho} \right) - k^{ED} C^{AUB} \quad (S14)$$

The linear biomass growth expression limits the total volume of biomass to 90% of the aqueous phase volume. At low biomass concentrations, such limits have negligible effects on biomass growth and substrate utilization because the biomass occupies a small volume of the total pore space.

When the biomass concentration begins to occupy a significant fraction of the pore volume, as might be expected near in-situ bioremediation injection wells, the key modeling assumption that biofilms in the pore space are thin and can be fully penetrated will likely be violated. The reduction (or near cessation) of biomass growth becomes less important than biofilm mass transport effects that are not considered in the model. Thus, through using the linear growth limitation expression, the model can only crudely approximate biological growth in grid blocks occupied by a substantial volume of biomass. At low biomass concentrations, the term has an insignificant effect.

The total biomass in the aquifer consists of the attached biomass and the unattached biomass is:

$$B^T = B^{AP} + \overline{B^A} \quad (S15)$$

where  $B^T$  is the total biomass,  $B^{AP}$  is the aqueous phase biomass, and  $\overline{B^A}$  is the attached biomass. The attached biomass is composed of the minimum biomass population ( $\overline{B^A}_{min}$ , which does not partition between the solid and the aqueous phase) and the biomass in equilibrium with the aqueous phase biomass:

$$\overline{B^A} = \overline{B^A}_{min} + \kappa B^{AP} \quad (S16)$$

Substituting the equilibrium relationship of equation (S15) into mass balance (S16) results in the following equilibrium concentration of aqueous phase biomass:

$$X = \frac{X_T - \overline{X}_{min}}{\kappa + 1} \quad (S17)$$

The attached biomass concentration is then calculated from equation (S16). The  $\kappa$  of infinity would mean that all of the biomass is attached, while the  $\kappa$  of 0 would mean that all of the biomass, except  $\overline{B^A}_{\min}$ , would exist in the aqueous phase.

#### Substrate competition

When two substrates (substrates 1 and 2) compete for the same enzyme, it reduces the rate of biodegradation. The half-saturation coefficient of each substrate in Monod term is suggested, and thus, the Monod terms for the two substrates would become (Chang & Alvarez-Cohen, 1995):

Substrate 1:

$$\frac{C_{S1}}{k_{S1}^{HS} \left(1 + \frac{C_{S2}}{k_{S2}^{HS}}\right) + C_{S1}} \quad (S18)$$

Substrate 2:

$$\frac{C_{S2}}{k_{S2}^{HS} \left(1 + \frac{C_{S1}}{k_{S1}^{HS}}\right) + C_{S2}} \quad (S19)$$

where  $C_{S1}$ ,  $C_{S2}$  are concentrations of substrates 1 and 2, respectively ( $\text{ML}^{-3}$ );  $k_{S1}^{HS}$ ,  $k_{S2}^{HS}$  are half-saturation coefficients of substrates 1 and 2, respectively ( $\text{ML}^{-3}$ ).

#### Inhibition

Inhibition effects can be addressed through multiplying the substrate biodegradation rate term by an inhibition factor (Chang & Alvarez-Cohen, 2010):

$$\left( \frac{I_{ih}}{I_{ih} + C_{ih}} \right) \quad (S20)$$

where  $I_{ih}$  is an experimentally determined inhibition constant for species  $ih$ . Inhibition can be used to simulate the sequential use of electron acceptors or the reduction of biodegradation rates due to the presence of a toxic or inhibitory compound. The term for substrate utilization and biomass growth can be calibrated by using one inhibition factor for each inhibiting substance.

#### Aerobic cometabolism

To describe the loss of cometabolite and attached biomass growth in aerobic cometabolic reactions, the following equations can be used (in the case of no mass transfer resistance, no inhibition, and no substrate competition) (de Blanc, 1998):

$$\frac{dC^{APC}}{dt} = -R^{SCB} C^{AUB} \left( \frac{C^{APC}}{k^{CHS} + C^{APC}} \right) \left( \frac{C^{APE}}{k^{EHS} + C^{APE}} \right) \left( \frac{C^{RP}}{K^{RHS} + C^{RP}} \right) \quad (S21)$$

$$\begin{aligned} \frac{dC^{\text{AUB}}}{dt} = & \frac{\mu_{\text{max}} C^{\text{AUB}}}{k^{\text{Y}}} \left( \frac{C^{\text{APS}}}{k^{\text{SHS}} + C^{\text{APS}}} \right) \left( \frac{C^{\text{APE}}}{k^{\text{EHS}} + C^{\text{APE}}} \right) \left( \frac{C^{\text{RP}}}{K^{\text{RHS}} + C^{\text{RP}}} \right) \left[ \frac{0.9(1 - C^{\text{AUB}})}{\rho} \right] \\ & - \frac{R^{\text{SCB}} C^{\text{AUB}}}{k^{\text{TC}}} \left( \frac{C^{\text{APC}}}{k^{\text{CHS}} + C^{\text{APC}}} \right) \left( \frac{C^{\text{APE}}}{k^{\text{EHS}} + C^{\text{APE}}} \right) \left( \frac{C^{\text{RP}}}{K^{\text{RHS}} + C^{\text{RP}}} \right) - k^{\text{ED}} C^{\text{AUB}} \end{aligned} \quad (\text{S22})$$

where  $R^{\text{SCB}}$  is maximum specific cometabolite biodegradation rate ( $\text{ML}^{-3}\text{T}^{-1}$ );  $C^{\text{APC}}$  is aqueous phase cometabolite concentration ( $\text{ML}^{-3}$ );  $C^{\text{RP}}$  is reducing power (NAD(P)H) concentration within the cells ( $\text{mMOL e}^-/\text{mass biomass}$ );  $K^{\text{RHS}}$  is NAD(P)H half-saturation constant ( $\text{mMOL e}^-/\text{mass biomass}$ );  $k^{\text{CHS}}$  is cometabolite half-saturation coefficient ( $\text{ML}^{-3}$ );  $\mu_{\text{max}}$  is maximum specific growth rate on growth substrate ( $\text{T}^{-1}$ ); and  $k^{\text{TC}}$  is transformation capacity (mass cells deactivated/mass cometabolite biodegraded).

**Text S5. Procedures for solving the coupled flow and transport problem**

The solution procedures are as follows:

- Step 1. Solve the pressure equation implicitly using a Jacobi conjugate gradient solver to yield water phase pressure in all grid blocks;
- Step 2. Capillary pressures from previous time step are used to determine the pressure of other phases in each grid block once the water phase pressure is known;
- Step 3. The Darcy's law is used to determine the phase velocities;
- Step 4. Mass conservation equations are solved explicitly to yield concentration of each component in each grid block;
- Step 5. Phase concentrations and saturations are determined through flash calculations;
- Step 6. The concentration of the components calculated by the pollutant migration model was used as the initial condition of the biodegradation model to obtain the pollutant degradation rate for this time step.
- Step 7. New capillary pressures are determined from the new saturations;
- Step 8. Repeat the procedures for each time step until simulation ends.

## Text S6. Stepwise cluster analysis (SCA)

In the stepwise-cluster analysis, the solutions of the numerical model (benzene concentrations at concerned locations) are considered as dependent variables; the operating conditions are independent variables. If the developed simulation model is run under  $n$  scenarios of system conditions, there will then be  $n$  sets of such independent and dependent variables (e.g., if the model is run 50 times under various system conditions, then  $n = 50$ ). Assume that there are  $m$  independent variables [e.g., four process control variables, denoted as  $x = (x_1, x_2, \dots, x_m)$ , where  $m = 4$ ], and  $p$  dependent variables [e.g., benzene concentrations at six concerned locations, denoted as  $y = (y_1, y_2, \dots, y_p)$ , where  $p = 6$ ]. Thus, all data can be given by matrixes  $X = (x_{tr})_{n \times m}$  and  $Y = (y_{tr})_{n \times p}$ , where  $r = 1, 2, \dots, m$ , and  $i = 1, 2, \dots, p$ .

The first step is to determine the clustering principles for the patterns. In SCA, patterns of responses will be cut or merged into a number of sets, and explanatory variables will be the references in judging which pattern set in the parent set should enter. After completion of cutting and merging processes, cluster trees could be produced and further used for predicting responses according to new explanatory values. The essence of this method is, based on a given criteria, to cut one pattern set of responses into two, and to merge two sets into one, step by step, in order to classify sets and sieve variables. Let cluster  $h$ , which contains  $n_h$  patterns, be cut into two sub-clusters  $e$  and  $f$ , containing  $n_e$  and  $n_f$  patterns, respectively (i.e.,  $n_e + n_f = n_h$ ). According to Wilks' likelihood-ratio criterion, if the cutting point is optimal, the value of Wilks'  $\Lambda$  ( $\Lambda = |W|/|T|$ ) should be minimum (Wilks, 1960; 1962; 1963; Kennedy and Gentle, 1981), where  $T$  and  $W$  are total-sample sum of the squares and cross products (SSCP) matrix  $\{t_{ij}\}$  and within-groups SSCP matrix  $\{w_{ij}\}$ , respectively, and  $T$  and  $W$  mean determinants of matrixes  $\{t_{ij}\}$  and  $\{w_{ij}\}$ , respectively. When the  $\Lambda$  value is very large, clusters  $e$  and  $f$  cannot be cut, but must be merged into greater cluster  $h$ . By Rao's  $F$ -approximation ( $R$ -Statistic), we have:

$$R = \frac{1-\Lambda^{1/S}}{\Lambda^{1/S}} \cdot \frac{Z \cdot S - P \cdot (K-1)/2 + 1}{P \cdot (K-1)} \quad (S23)$$

$$Z = n_h - 1 - (P + K)/2 \quad (S24)$$

$$S = \frac{P^2 \cdot (K-1)^2 - 4}{P^2 + (K-1)^2 - 5} \quad (S25)$$

where statistic  $R$  is distributed approximately as an  $F$ -value with  $v_1 = P \cdot (K - 1)$  and  $v_2 = P \cdot (K - 1)/2 + 1$  degrees of freedom,  $K$  is number of groups, and  $P$  is number of responses. The  $R$  - statistics will reduce to an exact  $F$ -value when  $P = 1$  or 2, or when  $K = 2$  or 3. Since the number of groups is two ( $K = 2$  for system operating conditions and benzene concentrations at concerned locations) in this study, an exact  $F$ -test is possible based on Wilks'  $\Lambda$  criterion. Thus, we have:

$$F(P, n_h - P - 1) = \frac{1-\Lambda}{\Lambda} \cdot \frac{n_h - P - 1}{P} \quad (S26)$$

Therefore, the criteria of cutting and merging clusters become to make a number of  $F$ -tests (Rao, 1952).

The second step is to test optimal cutting points, for which  $n_h$  patterns in cluster  $h$  are

sequenced according to the value of  $x_{r,k}^{(h)}$  in  $\{x_r\}$ , i.e.,  $x_{r,1^r}^{(h)} \leq x_{r,2^r}^{(h)} \leq \dots \leq x_{r,n_h^r}^{(h)}$ . Then the total-pattern SSCP matrix and within-groups SSCP matrix of responses  $y$  are calculated based on the sequence statistic  $\{K^r\}$ :

$$b_{ij}(K^r, n_h^r) = \frac{n_h^r K^r \cdot \{[B_i^{(h)}(K^r) - B_i^{(h)}(n_h^r)] \cdot [B_j^{(h)}(K^r) - B_j^{(h)}(n_h^r)]\}}{n_h^r - K^r} \quad (S27)$$

$$t_{ij}(n_h^r) = A_{ij}^{(h)}(n_h^r) - n_h^r B_i^{(h)}(n_h^r) B_j^{(h)}(n_h^r) \quad (S28)$$

$$w_{ij}(K^r, n_h^r) = t_{ij}(n_h^r) - b_{ij}(K^r, n_h^r) \quad (S29)$$

where:

$$B_{i \text{ or } j}^{(h)}(u) = \frac{1}{u} \sum_{k=1}^u y_{i \text{ or } j, k}^{(h)} \quad (S30)$$

$$A_{ij}^{(h)}(u) = \sum_{k=1}^u y_{i, k}^{(h)} y_{j, k}^{(h)} \quad (S31)$$

$$k^r = 1^r, 2^r, \dots, (n_h^r - 1), \forall r$$

$$i, j = 1, 2, \dots, p, \text{ and } r = 1, 2, \dots, m$$

For each  $x_r$ , a cutting point  $k^{*r}$  is derived, which satisfies:

$$\Lambda(k^{*r}, n_h^r) = \min_{k^r=1^r}^{(n_h^r-1)} \{\Lambda(k^r, n_h^r)\} \quad (S32)$$

For each explanatory variable, the index of response that will be used for cutting judgments (denoted as  $r^*$ ) is derived, which satisfies:

$$\Lambda(k^{*r^*}, n_h^r) = \min_{r=1}^m \{\Lambda(k^r, n_h^r)\} \quad (S33)$$

Thus, the optimal cutting point of cluster  $h$  is  $k^{*r^*}$ , and the relevant value of explanatory variable (i.e., the reference for new pattern prediction) is  $x_{r^*, k^{*r^*}}^{(h)}$ . Then a  $F$ -test can be undertaken.

If

$$F(P', n_h^{r^*} - P' - 1) = \frac{1 - \Lambda(k^{*r^*}, n_h^{r^*})}{\Lambda(k^{*r^*}, n_h^{r^*})} \frac{n_h^{r^*} - P'}{P'} \geq F_1 \quad (S34)$$

is satisfied, cluster  $h$  can be cut into two sub-clusters according to the distribution of  $x_{r^*}$ : (a) data in explanatory sets with  $k^{r^*} \leq k^{*r^*}$  are allocated into sub-cluster  $e$  ( $< f$ ); (b) data in explanatory sets with  $k^{r^*} > k^{*r^*}$  are allocated into sub-cluster  $f$ , where  $P'$  is number of responses under consideration. Among explanatory variables,  $x_{r^*}$  is the most important one affecting the response. If equation (S12) is not satisfied, cluster  $h$  cannot be cut. Then the other clusters will be tested to decide whether to cut or not, i.e., to test  $h = 1, 2, \dots, H$  ( $H$  is total number of clusters at the current stage). When no cluster can be cut, the next step is to merge the clusters.

The third step is the mergence of clusters. To test the mergence of clusters  $e$  and  $f$  for existing clusters, the total-sample SSCP matrix and within-groups SSCP matrix should be calculated firstly:



$$t_{ij}(n_e, n_f) = A_{ij}^{(e)}(n_e) + A_{ij}^{(f)}(n_f) - \left[ n_e B_i^{(e)}(n_e) + n_f B_i^{(f)}(n_f) \right] \cdot \left[ n_e B_j^{(e)}(n_e) + n_f B_j^{(f)}(n_f) \right] / (n_e + n_f) \quad (S35)$$

$$b_{ij}(n_e, n_f) = \frac{n_e n_f [B_i^{(e)}(n_e) - B_i^{(f)}(n_f)] [B_j^{(e)}(n_e) - B_j^{(f)}(n_f)]}{(n_e + n_f)} \quad (S36)$$

$$w_{ij}(n_e, n_f) = t_{ij}(n_e, n_f) - b_{ij}(n_e, n_f) \quad (S37)$$

where  $A_{ij}$  and  $B_i$  or  $B_j$  have the same formulation as equations (S30) and (S31);  $i, j = 1, 2, \dots, p$ . Then a  $F$ -test can be undertaken. If

$$F(P', n_e + n_f - P' - 1) = \frac{1 - \Lambda(n_e, n_f)}{\Lambda(n_e, n_f)} \frac{(n_e + n_f)^{P' - 1}}{P'} < F_2 \quad (S38)$$

is satisfied, clusters  $e$  and  $f$  can be merged into a new cluster  $h$ . Otherwise, it should be similarly tested whether other clusters can be merged for  $e = 1, 2, \dots, (H-1)$  and  $f = 2, 3, \dots, H$ .

The final step is the prediction of the response according to new explanatory variables. After all calculations and tests have been completed (i.e., all hypotheses of further cutting or merge are rejected), a cluster tree can be derived for each response. Each cutting point, which leads to two branches, corresponds to the value  $(x_{r^*, k^*}^{(h)})$  of an explanatory variable. When a new pattern set of explanatory variables  $\{x_r\}$  is examined, its  $x_{r^*}$  value can be compared with  $x_{r^*, k^*}^{(h)}$  at the cutting point, and classified into relevant branches. Step-by-step, the pattern will finally enter a tip cluster which cannot be either cut or merged further. The criterion to classify a new sample to relevant branches is that, (a) sample data with  $x_{r^*} \leq x_{r^*, k^*}^{(h)}$  are merged into cluster  $e$  ( $< f$ ) and (b) sample data with  $x_{r^*} > x_{r^*, k^*}^{(h)}$  are merged into cluster  $f$ . Let  $e'$  be the tip cluster where the new sample enters. Then the predicted dependent variable  $\{y_i\}$  is:

$$y_i = y_i^{(e')} \pm R_i^{(e')} \quad (S39)$$

where  $y_i^{(e')}$  is mean of dependent variable  $i$  in sub-cluster  $e'$ , and  $R_i^{(e')}$  is radius of  $y_i$  in cluster  $e'$ :

$$y_i^{(e')} = \left\{ \max_{k=1}^{n_{e'}} (y_{i,k}^{(e')}) + \min_{k=1}^{n_{e'}} (y_{i,k}^{(e')}) \right\} / 2, \forall i \quad (S40)$$

$$R_i^{(e')} = \left\{ \max_{k=1}^{n_{e'}} (y_{i,k}^{(e')}) - \min_{k=1}^{n_{e'}} (y_{i,k}^{(e')}) \right\} / 2, \forall i \quad (S41)$$

## Text S7. Filtering Process Model

After the clustering process, a number of leaf clusters are produced. Each leaf cluster contains a group of modeling outputs with similar statistical attributes; these modeling outputs provide an output value range for the leaf cluster. The purpose of filtering is to calculate an optimal estimate for each leaf cluster; this estimate can be used as an optimal output value for the leaf cluster. The set of leaf clusters for all well patterns thus can be regarded as all possible results for the remediation design.

Among various filtering methods, the well-known Kalman filter has been recognized as a powerful tool in supporting estimations of past, present, and future states. In this study, a filtering process model based on the Kalman filter method was developed to calculate the optimal estimate for each leaf cluster.

Generally, the Kalman filter addresses the problem of estimating the state of a discrete-time controlled process,  $z$  ( $z \in R^f$ ), that is governed by the following linear stochastic difference equation:

$$z_k = Az_{k-1} + Bu_{k-1} + w_{k-1} \quad (S42)$$

with a measurement ( $q \in R^g$ ) as follows:

$$q_k = Hz_k + v_k \quad (S43)$$

where  $u_k$  is the optional control input ( $u \in R^l$ );  $w_k$  and  $v_k$  represent the process and measurement noise (random variables), respectively. They are assumed to be independent (of each other), white, and with normal probability distributions  $p(w) \sim N(0, Q^{PNC})$  and  $p(v) \sim N(0, R^{MNC})$ , respectively. A white noise process is defined as a random process of random variables that are uncorrelated, have mean zero, and a finite variance. The process noise covariance  $Q^{PNC}$  and measurement noise covariance  $R^{MNC}$  matrices are assumed to be constant.

In equation (S42), the  $f \times f$  matrix  $A$  relates the state at the previous time step ( $k-1$ ) to the state at the current step  $k$ , in the absence of a process noise. The  $f \times l$  matrix  $B$  relates the optional control input  $u$  to the state  $z$ . The  $g \times f$  matrix  $H$  in equation (S21) relates the state to the measurement  $q$ . Matrices  $A$  and  $H$  are assumed to be constants.

It is defined that  $\hat{z}_{\bar{k}}$  ( $\hat{z}_{\bar{k}} \in R^f$ ) is a priori-state estimate at step  $k$  given knowledge of the process prior to step  $k$ , and  $\hat{z}_k$  ( $\hat{z}_k \in R^f$ ) to be a posteriori-state estimate at step  $k$  given measurement  $q_k$ . It is then defined a priori-estimate error and a posteriori-estimate error as  $e_{\bar{k}} \equiv z_k - \hat{z}_{\bar{k}}$  and  $e_k \equiv z_k - \hat{z}_k$ , respectively. Thus, the priori-estimate error covariance can be written as  $P_{\bar{k}} = E[e_{\bar{k}}e_{\bar{k}}^T]$  and the posteriori-estimate error covariance as  $P_k = E[e_k e_k^T]$ .

The posteriori-state estimate ( $\hat{z}_k$ ) can be calculated as:

$$\hat{z}_k = \hat{z}_{\bar{k}} + K(q_k - H\hat{z}_{\bar{k}}) \quad (S44)$$

The difference between the actual measurement ( $q_k$ ) and the measurement prediction, ( $q_k - H\hat{z}_{\bar{k}}$ ), in equation (S44) is called the residual, which reflects the discrepancy between the predicted measurement and the actual measurement.

The  $f \times g$  matrix  $K$  in equation (S44) is Kalman gain, which is chosen to minimize the posteriori error covariance. The Kalman gain  $K_k$  can be given as follows (Maybeck, 1979; Jacobs, 1993):

$$K_k = P_{\bar{k}} H^T (H P_{\bar{k}} H^T + R)^{-1} \quad (\text{S45})$$

As the  $R^{MNC}$  approaches zero, the gain  $K$  weights the residual more heavily. Specifically,  $\lim_{R_k \rightarrow 0} K_k = H^{-1}$ . On the other hand, as the priori-estimate error covariance  $P_{\bar{k}}$  approaches zero, the gain  $K$  weights the residual less heavily. Specifically,  $\lim_{P_{\bar{k}} \rightarrow 0} K_k = 0$ .

The Kalman filter consists of time-update equations and measurement-update equations. The discrete time-update equations are written as:

$$\hat{z}_{\bar{k}} = A \hat{z}_{k-1} + B u_{k-1} \quad (\text{S46})$$

$$P_{\bar{k}} = A P_{k-1} A^T + Q^{PNC} \quad (\text{S47})$$

The time-update equations are responsible for projecting forward (in time) the current state and error covariance estimates to obtain the priori-estimates for the next time step. The discrete measurement-update equations are given as:

$$K_k = P_{\bar{k}} H^T (H P_{\bar{k}} H^T + R^{MNC})^{-1} \quad (\text{S48})$$

$$\hat{z}_k = \hat{z}_{\bar{k}} + K_k (q_k - H \hat{z}_{\bar{k}}) \quad (\text{S49})$$

$$P_k = (1 - K_k H) P_{\bar{k}} \quad (\text{S50})$$

The measurement-update equations are responsible for the feedback—i.e., for incorporating a new measurement into the priori-estimate to obtain an improved posteriori-estimate.

The operation of the filter is shown below. The first step is to compute the Kalman gain,  $K_k$ . The next step is to actually measure the process to obtain  $q_k$  and then to generate a posteriori-state estimate by incorporating the measurement as in equation (S49). The final step is to obtain a posteriori error covariance estimate via equation (S50). After each iteration of time update and measurement update, the process is repeated with the previous posteriori-estimates used to predict the new priori-estimates. This recursive nature is one of the very appealing features of the Kalman filter. For example, compared with the implementation of a Wiener filter, which operates on all the data directly for each estimate, the implementation of the Kalman filter is much more feasible.

In this study, the modeling outputs (samples) in each leaf cluster can be regarded as measurements. For any leaf cluster, the time update equations were written as:

$$\hat{z}_{\bar{k}} = \hat{z}_{k-1} \quad (\text{S51})$$

$$P_{\bar{k}} = P_{k-1} + Q^{PNC} \quad (\text{S52})$$

where  $A=1$  (the state did not change from step to step), and  $u=0$  (there was no control input). The measurement update equations were given as:

$$K_k = P_{\bar{k}} (P_{\bar{k}} + R^{MNC})^{-1} \quad (\text{S53})$$

$$\hat{z}_k = \hat{z}_{\bar{k}} + K_k (q_k - \hat{z}_{\bar{k}}) \quad (\text{S54})$$

$$P_k = (1 - K_k) P_{\bar{k}} \quad (\text{S55})$$

where  $H=1$  (the noisy measurement is of the state directly);  $k$  denotes the number of samples (modeling outputs) in each leaf cluster.

After the clustering and filtering, an optimal estimate can be obtained for each leaf cluster. A new sample can be grouped into a corresponding leaf cluster by comparing the values of  $x_{tr}$  with those of  $x_{r^*}^{\alpha}(h^*)$ . The corresponding output variable can be predicted as  $y_i = \hat{z}_{k,i}, \forall i$ .

### Text S8. Nonlinear Optimization model of the FCI optimizer

The Nonlinear Optimization model of the FCI optimizer can be formulated as follows (to identify the optimum control conditions):

$$\text{Min } Z = \sum_{i=1}^I U_i^{In} + \sum_{j=1}^J U_j^{Ex} \quad (\text{S56})$$

subject to:

$$X_{kt}(U_i^{In}, U_j^{Ex}) \leq X_{max} \text{ for all } k=1,2, \dots, K \quad (\text{S57})$$

$$0 \leq U_i^{In} \leq U_{i,max}^{In} \quad (\text{S58})$$

$$0 \leq U_j^{Ex} \leq U_{j,max}^{Ex} \quad (\text{S59})$$

$$\sum_{i=1}^I U_i^{In} = \sum_{j=1}^J U_j^{Ex} \quad (\text{S60})$$

where  $Z$  is the total pumping rate for all injection and extraction wells;  $U_i^{In}$  and  $U_j^{Ex}$  are pumping rates for the  $i$ th injection well and the  $j$ th extraction well after a period of remediation;  $U_{i,max}^{In}$  and  $U_{j,max}^{Ex}$  are maximum pumping rates for the  $i$ th injection well and the  $j$ th extraction well;  $X_{max}$  is environmental standard;  $I, J, K$  are numbers of injection well, extraction well, and monitoring well, respectively;  $X_{kt}$  is predicted benzene concentration at  $t$ . Constraint (S60) indicates that all the extracted water will be injected into the aquifer. This constraint is emphasized to ensure such a stable hydraulic gradient that the groundwater can flow directed toward the plume interior.

### Text S9. Nonlinear Optimization model of the DPC system

The Nonlinear Optimization model of the DPC system can be formulated as follows (to identify the optimum control conditions):

$$\text{Min } Z = w_1(X)(S(X) - H)^2 + w_2(U)U \quad (\text{S61})$$

subject to:

$$S(X) = (X - X^0) / X^0 \quad (\text{S62})$$

$$X = F(U) \quad (\text{S63})$$

$$U_L \leq U \leq U_U \quad (\text{S64})$$

where  $Z$  is the optimization objective, representing the system cost;  $U_L$  and  $U_U$  are the lower and upper bounds of  $U$ , respectively;  $w_1$  and  $w_2$  are the weights to reflect different priorities for the remediation efficiency and cost. In this optimization model,  $S(X)$  is within the range of 0 to 1; therefore, the injection and extraction rates ( $U$ ) are normalized to fit it.  $H$  is a constant greater than or equal to 1 which is the highest contaminant removal rate. In this optimization model, a pseudo-equation  $X=F(U)$  is used to describe the relationship between  $X$  and  $U$ .

### Text S10. Nonlinear Optimization model of the SADPC system

The Nonlinear Optimization model of the SADPC system can be formulated as follows (to update the optimum control conditions):

$$\text{Min } J = \left\{ \sum_{i=1}^p \omega_i(X) (X_r(t+i) - X_p(t+i))^2 + \sum_{i=1}^p \omega_i(U) U(t+i-1) \right\} \quad (\text{S65})$$

subject to:

$$X = F(U) \quad (\text{S66})$$

$$0 \leq X_r(t+i) \leq X_{max} \quad (\text{S67})$$

$$0 \leq X_p(t+i) \leq X_{max} \quad (\text{S68})$$

$$U_L \leq U \leq U_U \quad (\text{S69})$$

where  $J$  is the optimization objective, representing the system cost;  $P$  is the prediction horizon;  $w_i(X)$  and  $w_i(u)$  are the weights to reflect different priorities for the remediation efficiency and cost.  $X_r(t+i)$  and  $X_p(t+i)$  are setpoint and predicted value, respectively;  $X_{max}$  is environmental standard.  $U$  is the operating condition;  $U_L$  and  $U_U$  are the lower and upper bounds of  $U$ , respectively. In this optimization model, a pseudo-equation  $X = F(U)$  is used to describe the relationship between  $X$  and  $U$ .

## **Text S11. Genetic algorithms (GA)**

GAs are heuristic search procedures based on the mechanisms of genetics and Darwin's natural selection principles, combining an artificial survival of the fittest with genetic operators abstracted from nature (Holland, [1975](#)).

An initial random population of genomes within the search space is generated. Each genome represents a possible solution to the search/optimization problem and is represented by a string of values (genes), one per each search variable. Survival of the fittest is accomplished by evaluating each genome's fitness through an appropriate objective function and a biased random selection procedure of individuals for "reproduction", where higher rated genomes are more likely to be selected. Generation of a new population is achieved by means of crossover (partial exchange of information between pairs of strings) and mutation (a random change in a random location within the string). The fittest individuals are transferred unchanged to the next generation, an approach known as "elitism". Every new generation of genomes is expected to be more closely concentrated in the vicinity of the optimal solution. The process is repeated until a convergence criterion is met or a pre-set maximum number of generations reached. GA input parameters include: population size, number of generations, range limits of each gene, crossover and mutation rates and a fitness function for genome evaluation.

In this study, GA is used to solve the developed discrete and nonlinear model. A set of parameters are needed to be predefined for guiding the genetic algorithm (Kuo et al., [2006](#); Matott et al., [2006](#); Stramer et al., [2010](#); Opher and Ostfeld, [2011](#); Liao et al., [2020](#)), including: (1) chromosome length LCHR which is the product of the number of decision variables ( $n$ ) and the length of a string ( $k$ ); (2) population size  $M$  which is usually within the range of 30 to 200; (3) crossover rate RCRO which is usually within the range of 0.6 to 0.95; (4) mutation rate RMUT which is usually within the range of 0.001 to 0.05; and (5) convergence criterion which is used to judge whether stop the search process. Normally the process is stopped after a predetermined generation number NG is reached or when there are no significant differences among the best solutions.



### Text S12. MPC control module procedure

The running procedures are as follows:

Step 1. Set the prediction time domain  $P$  and the weighting coefficients  $\omega_i$ ;

Step 2. Use the expected output sequence  $x_r(t)$  in the future, and the reference trajectory comes from the first-order exponential form fitting the actual output value of the DPC system;

Step 3. The control amount obtained in this sampling time period from the DPC system is brought into the biodegradation process to obtain the actual system output  $x(t)$ ;

Step 4. Use the DPC system to obtain the model output  $x_m(t)$  of the current sampling time period and the predicted output  $x_m(t+i)$  of the future time period, and obtain the system predicted output value  $x_p(t+i)$  after feedback correction;

Feedback correction: 
$$x_p(t+i) = x_m(t+i) + he(t) \quad (S70)$$

$$e(t) = x(t) - x_m(t) \quad (S71)$$

$h$  is the compensation coefficient;

Step 5. The optimization algorithm is used to solve the rolling optimization, and the optimal sequence  $U(t+i-1)$  is obtained;

Step 6. Apply the first control variable  $U(t)$  of the optimal sequence to the system, and then return to step 2.

### **Text S13. General procedure for developing a process control system for enhanced in situ biodegradation**

Step 1. A 3D pilot-scale model is designed for supporting the operation of enhanced in-situ biodegradation.

Step 2. After the occurrence of a hydrocarbon spill, an enhanced in-situ biodegradation process is to be undertaken. A subsurface LNAPLs biodegradation model is then developed to reflect the in-situ LNAPL biodegradation process.

Step 3. After calibration and verification, the interactions between contaminant concentrations and operating conditions are simulated through the subsurface model.

Step 4. Considering high complexities and computational requirements in incorporating numerical simulation model directly into optimization frameworks, coupled with the inability to obtain enough samples due to the high cost of sampling, a statistical relationship between remediation system performance and operating condition will be developed based on a large number of runs for the developed simulation model under various system conditions. Different scenarios of contamination situations and operating conditions are considered for the simulation. Under each contamination situation, the effects of various operation conditions on contaminant concentrations at concerned locations are examined.

Step 5. The stepwise cluster analysis method or the filtered clustering analysis method is used to develop to reflect the effects of variations of operating-condition on contaminant concentrations. Thus, a bridge between the subsurface model and the operating decision is established for further determining the desired operating conditions.

Step 6. Based on the established statistical relationships, a corresponding nonlinear discrete optimization model for groundwater control is established to determine optimal operating conditions corresponding to specific contamination situations. The GA technique is used to solve the developed optimization model.

Step 7. After the optimal operation conditions for each scenario are determined, the SI emulator is developed through the obtained knowledge base.

Step 8. A new nonlinear discrete optimization model is formulated by using the part that meets the expectation and its epitaxial as the setpoint curve for the contamination situations that do not meet the expectation in each scenario. Rolling optimization determines the optimal operating conditions corresponding to a specific contamination situation. The GA technique is used to solve the newly developed optimization model, and the optimal operating conditions of each scene are updated.

#### Text S14. Collinearity test of the independent variables

The presence of high collinearity in an FCI simulator implies that the conclusions of the analysis can be questioned. For example, the accuracy of estimations cannot be guaranteed due to high variances of the estimators. Thus, detection of collinearity should be a compulsory first step in every correlation analysis. Collinearity measures have been widely applied to examine if there are any co-relations among the independent variables. Variance Inflator Factor (VIF) is commonly used to evaluate the level of collinearity, which can be calculated as follows:

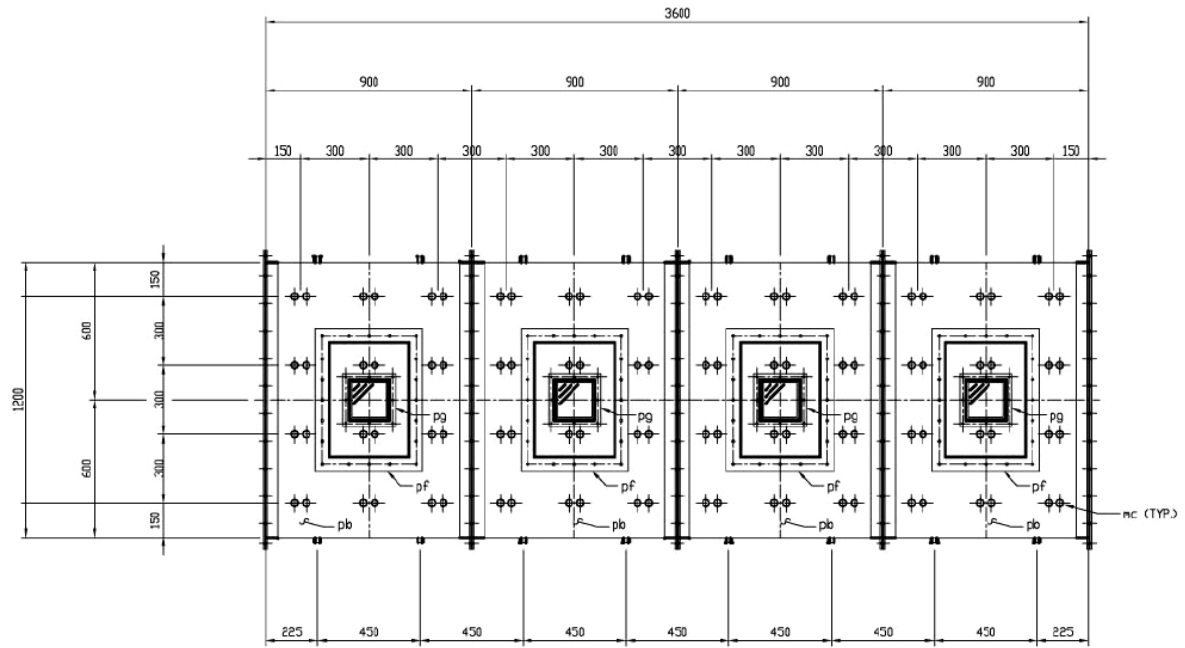
$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (\text{S72})$$

$$\mathbf{X} = (\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3, \mathbf{u}_4) \quad (\text{S73})$$

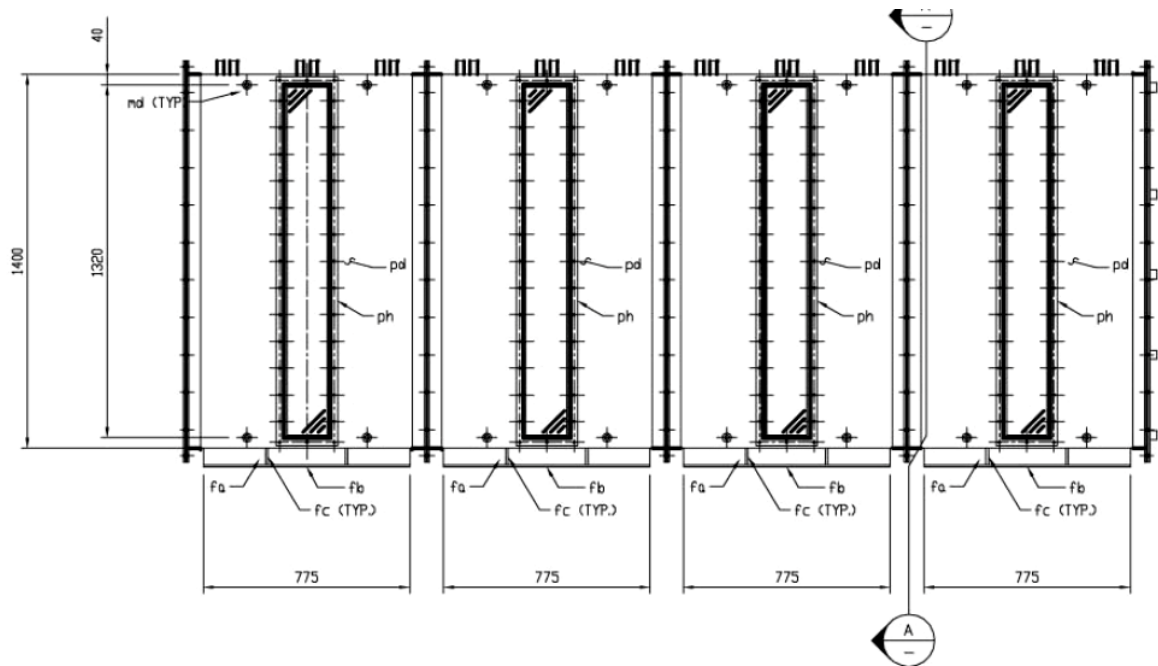
$$\text{VIF}_i = \frac{1}{1-R_i^2}, i = 1, 2, 3, 4 \quad (\text{S74})$$

where  $\boldsymbol{\varepsilon}$  represents the random disturbance and  $R_i$  is the negative correlation coefficient of the independent variable for the regression analysis of the remaining independent variables.

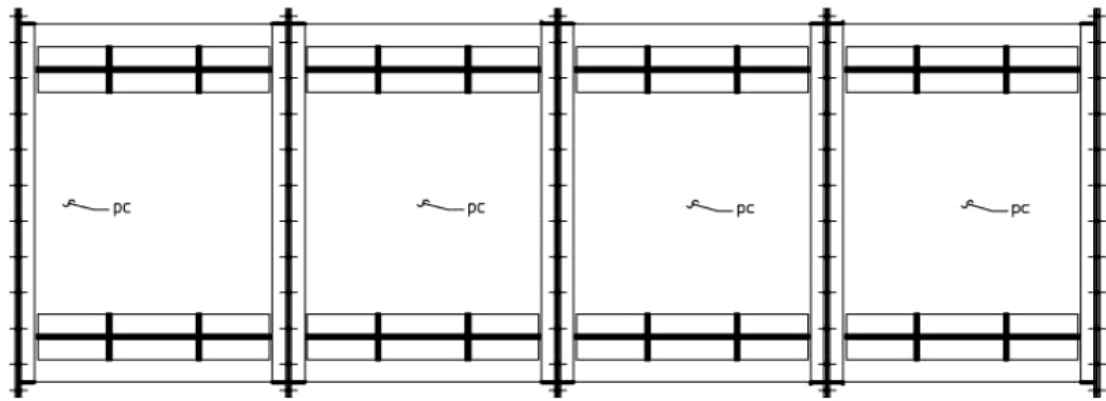
If the data matrix has no full column rank that can be considered “severe multicollinearity”, e.g., an independent variable can be expressed linearly by other independent variables. The closer the VIF value near 1, the lower the collinearity level is. The threshold value is usually 10. In this study, we selected groundwater injection rates of oxygen and nutrient in Well I and Well II ( $\mathbf{u}_1$  and  $\mathbf{u}_2$ ), and groundwater extraction rates in Well III and Well IV ( $\mathbf{u}_3$  and  $\mathbf{u}_4$ ) as the independent variables (called control variables in this paper) (Table S6). Results show that the corresponding VIF values of all the independent variables are much less than the threshold value (10), indicating that the variables are independent and do not have the multicollinearity (Table S8).



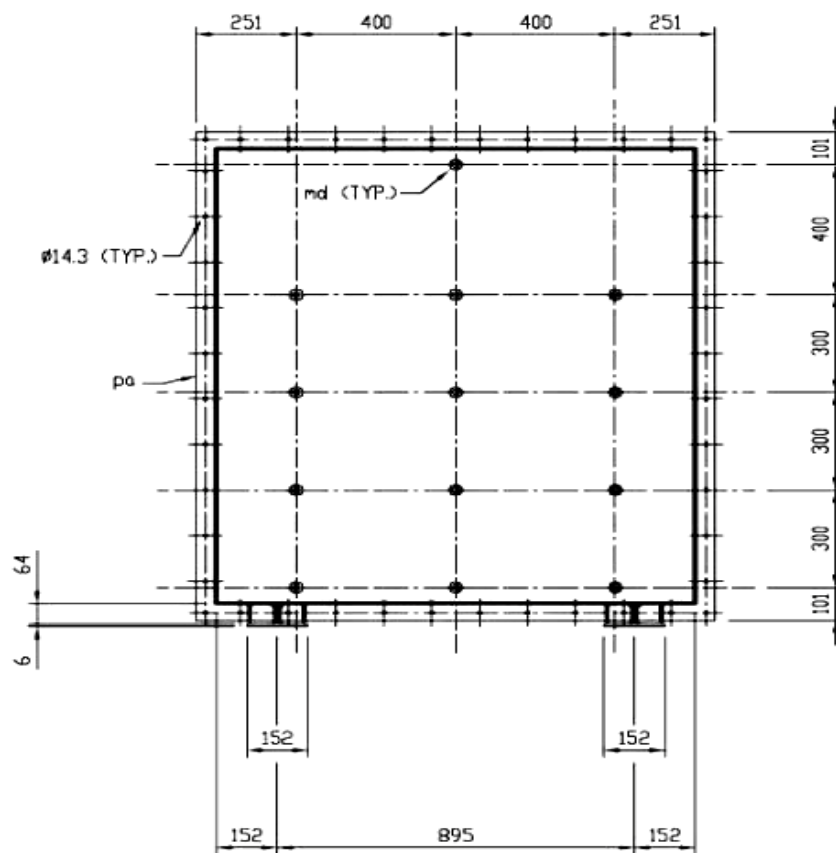
**Figure S1(a).** Plan view of the pilot scale system



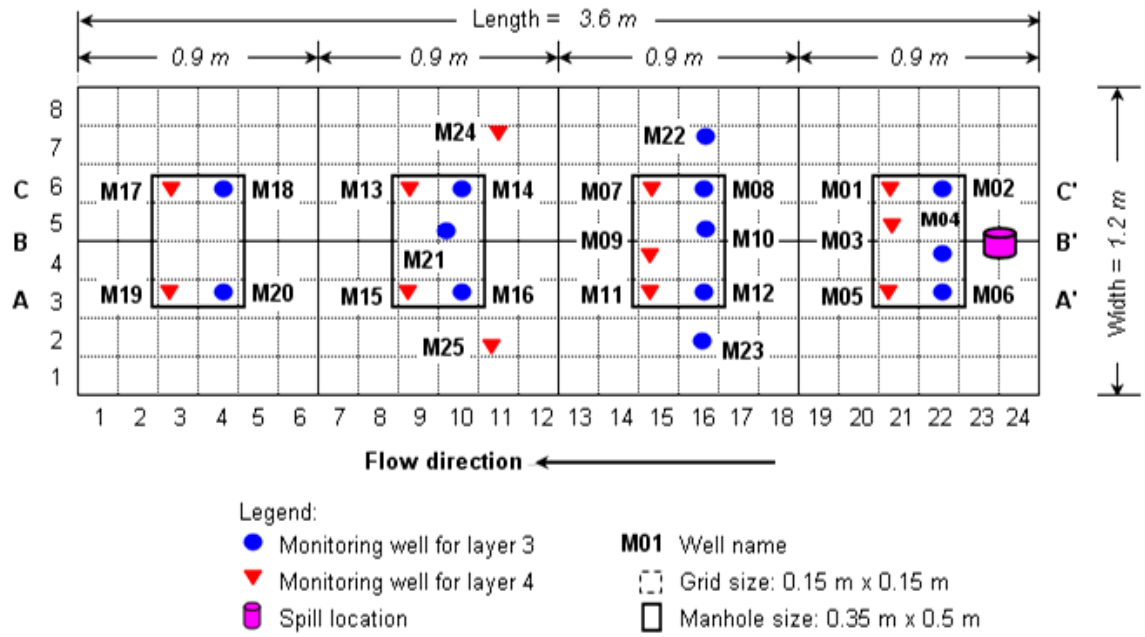
**Figure S1(b).** Front view of the pilot scale system



**Figure S1(c).** Bottom view of the pilot scale system

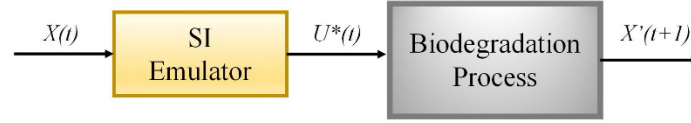


**Figure S1(d).** End elevation of the pilot scale system

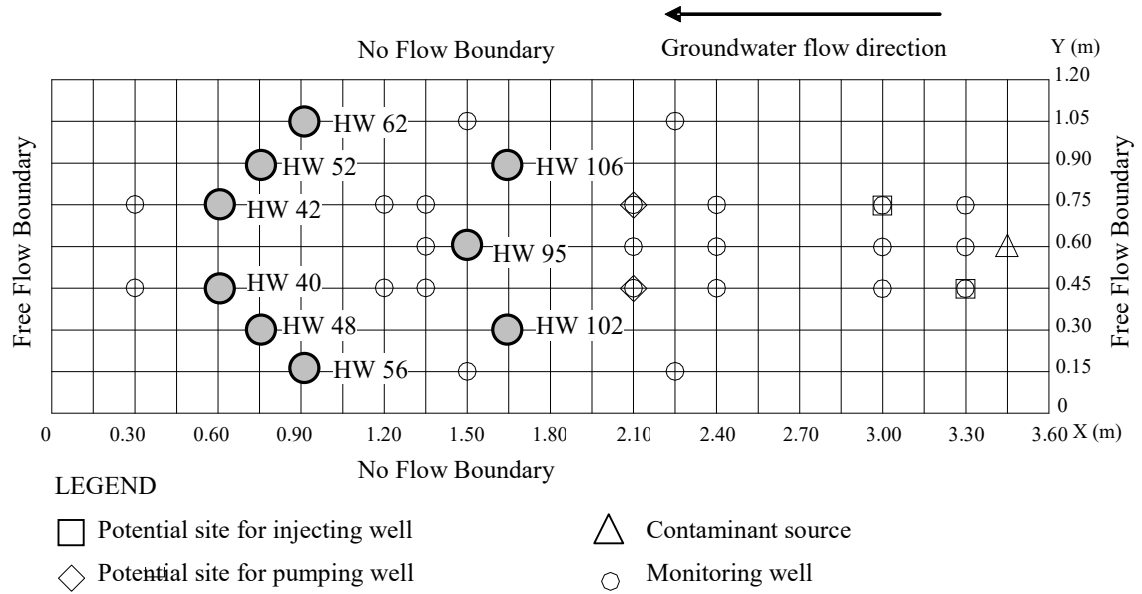


**Figure S1(e).** Well locations (plan view)

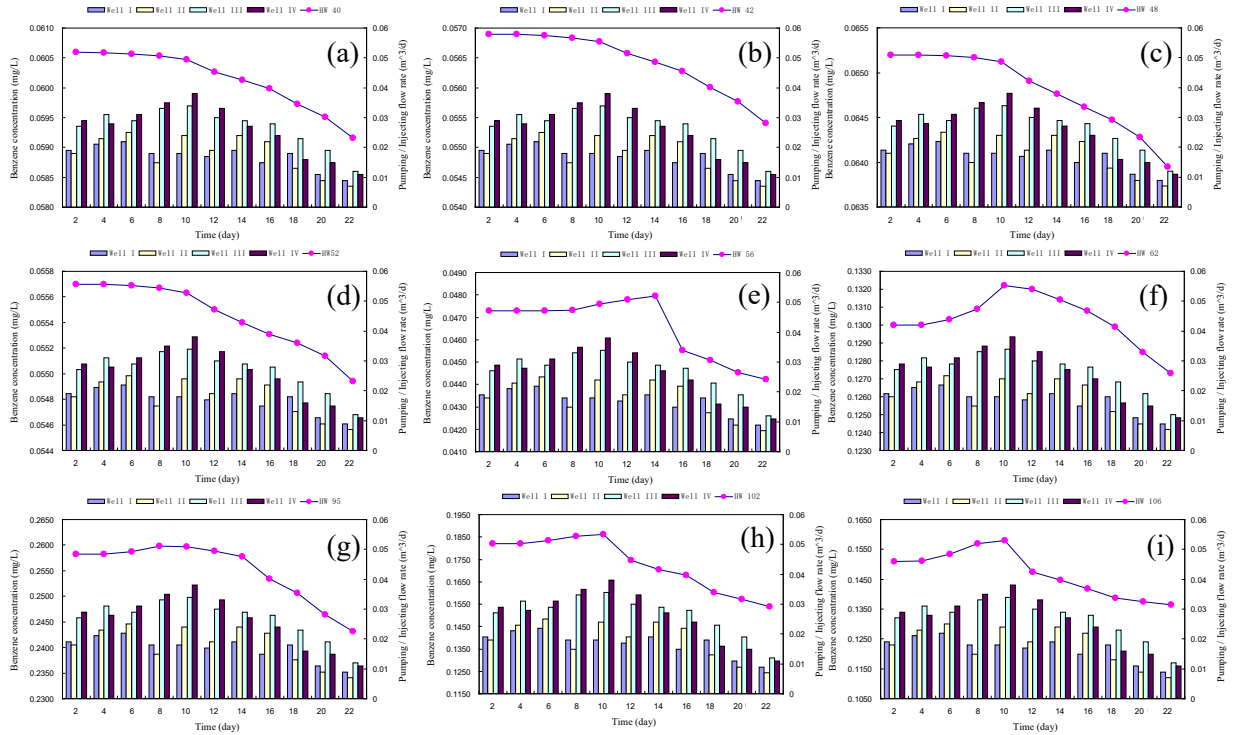




**Figure S3.** Framework of the SI Emulator

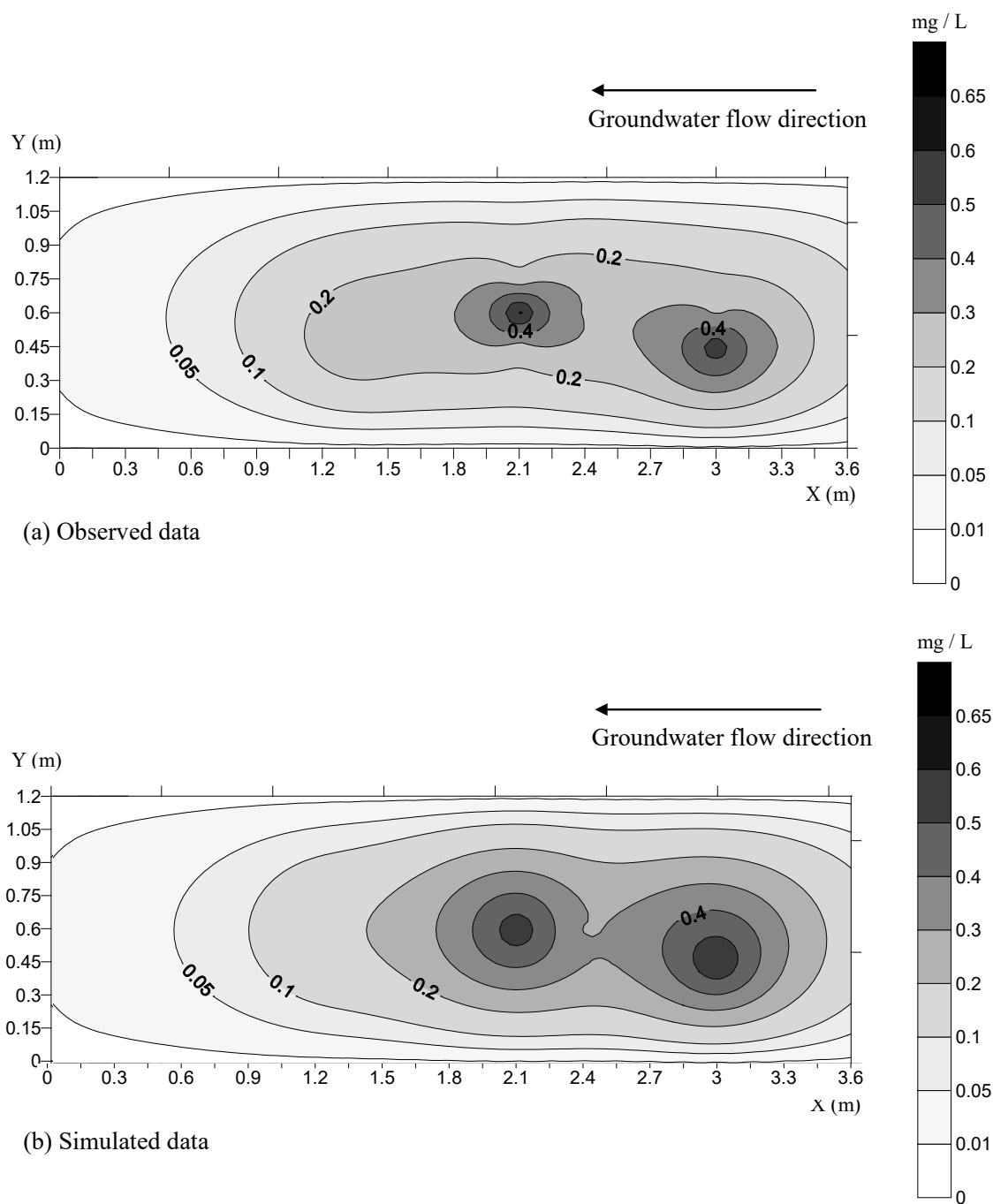


**Figure S4.** Locations of the hypothetical wells

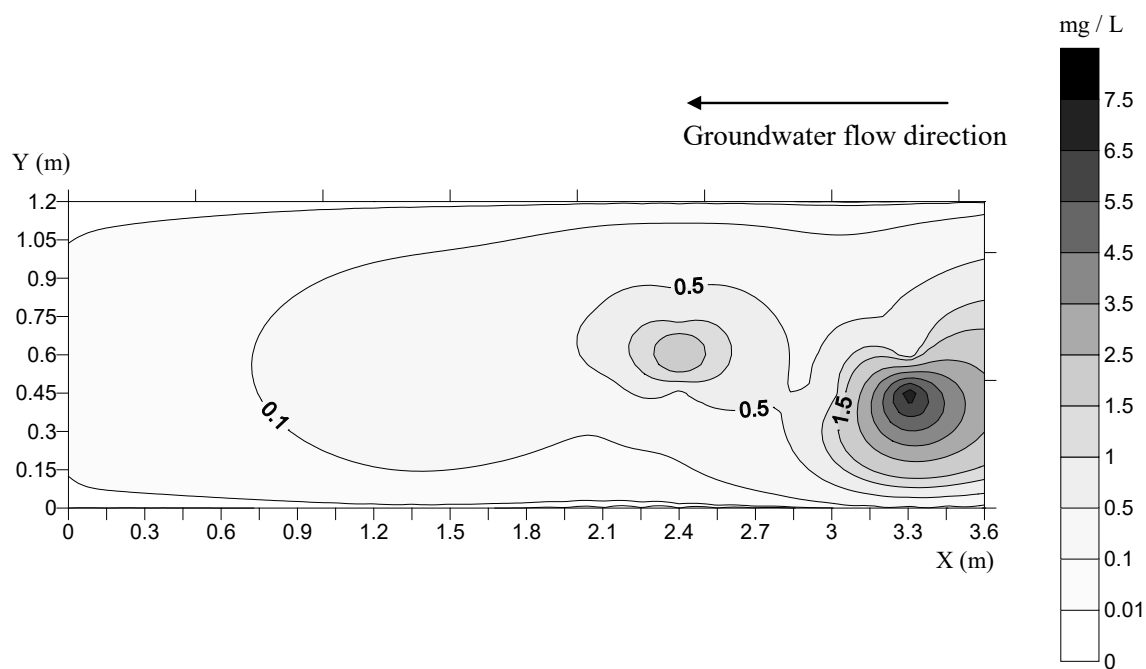




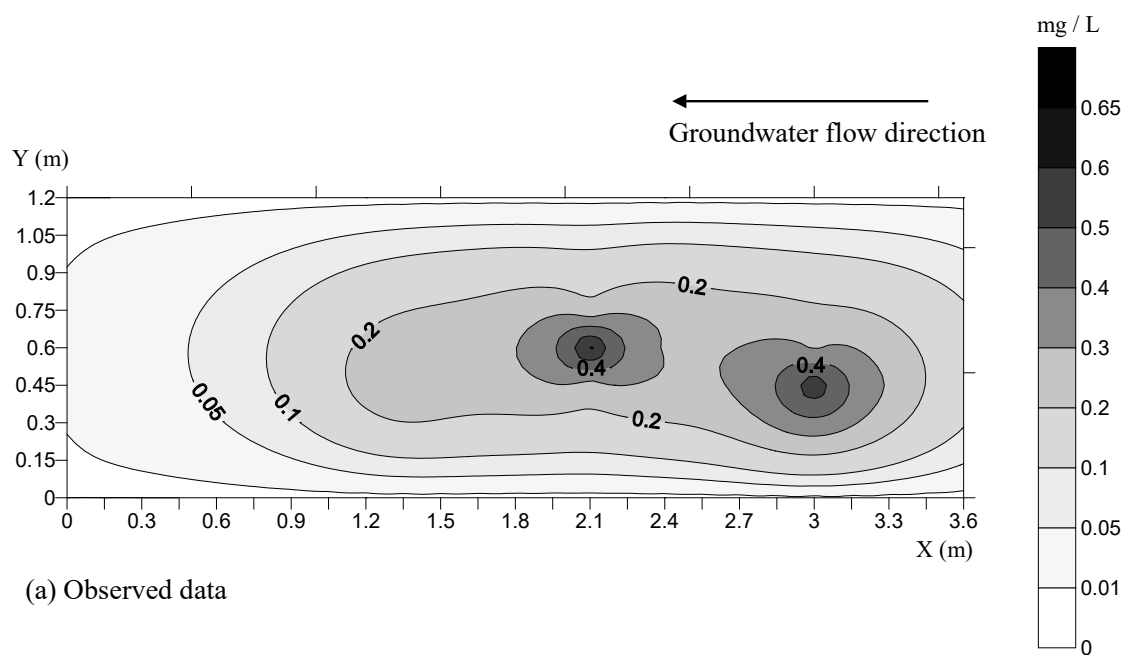
**Figure S5.** Benzene concentrations of the DPC system from Day 2 to Day 22, where Figs. (a) to (i) represents the concentrations at HW-40, HW-42, HW-48, HW-52, HW-56, HW-62, HW-95, HW-102, and HW-106



**Figure S6.** The concentration distribution of Benzene on Day 57 of the experiment

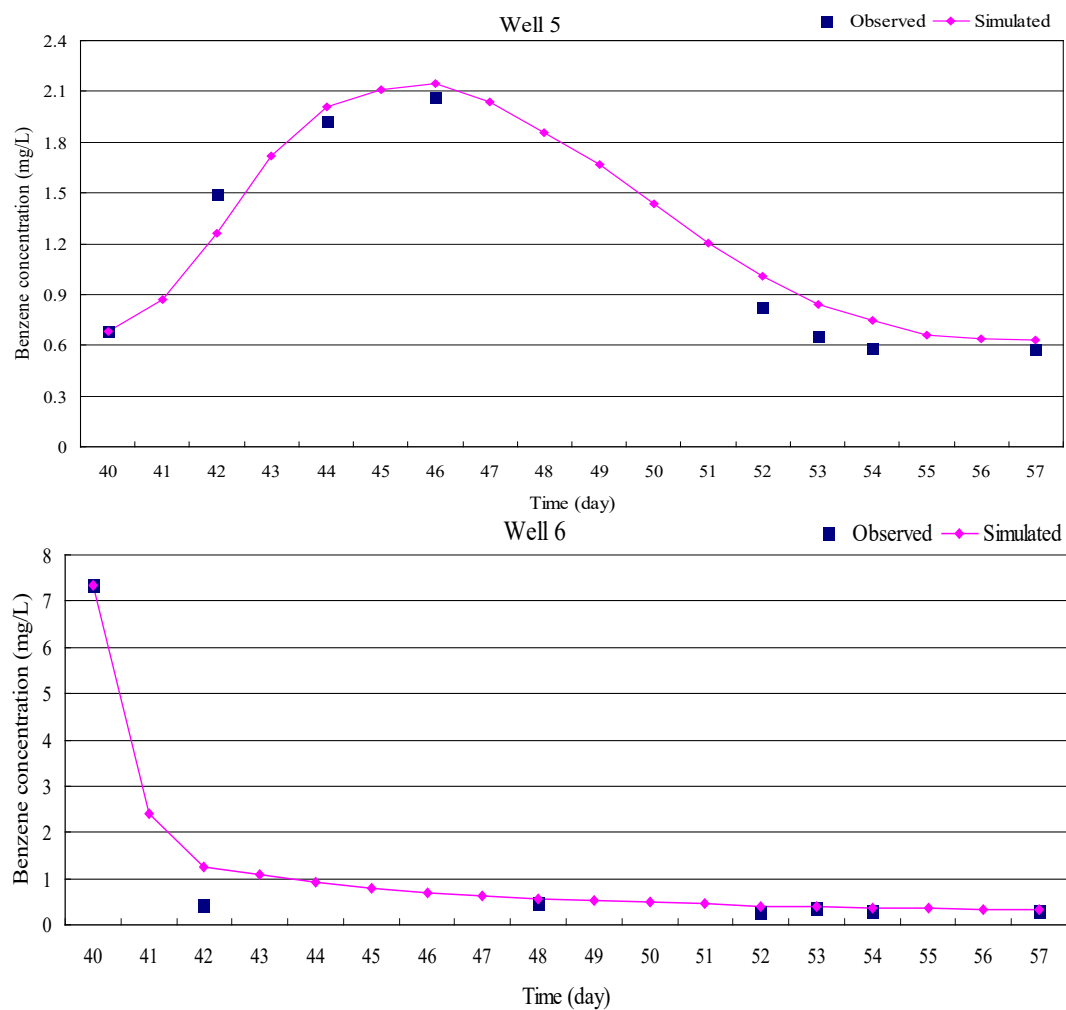


**Figure S7(a).** The concentration distribution of Benzene on Day 40



(a) Observed data

**Figure S7(b).** The concentration distribution of Benzene on Day 57



**Figure S8.** Verification results for well 5 and well 6

**Table S1.** Input parameters for contaminant transport simulation

Parameter	Value
<b>Flow and transport simulation parameters</b>	
Hydraulic conductivity of sand/ till/ clay	10 / 5 / 2.5 m/d
Permeability of sand/ till/ clay	1500/430/ 890 MD
Porosity of sand/ till/ clay	0.35 / 0.30 / 0.45
Longitudinal dispersivity of sand/ till/ clay	0.1 / 0.1 / 0.1 m
Transverse dispersivity of sand/ till/ clay	0.01 / 0.01 / 0.01 m
Van Genuchten's alpha of sand/ till/ clay	10 m <sup>-1</sup>
Van Genuchten's n of sand/ till/ clay	6.8
First-order reaction rate coefficient of benzene	0.21 /d
Endogenous decay coefficient	0.2544 /d
Residual water saturation	0.01
Water dynamic viscosity	1.0 cp
Water interfacial tension	45 Dynes/cm
Benzene density	0.713 g/cm <sup>3</sup>
Hydraulic gradient	0.03 m/m
Water partition coefficient of benzene	0.00203
Benzene solubility	1750 mg/L
Aquifer thickness	1.2 m
Time step	0.101 day
Maximum time step size	10 day
Tolerance for concentration change	0.001
<b>Enhanced biodegradation simulation parameters</b>	
Water injecting rate	20 L/d
NH <sub>4</sub> NO <sub>3</sub> nutrient injecting rate	1750 mg/L
NH <sub>4</sub> HPO <sub>4</sub> nutrient injecting rate	1100 mg/L
Heterotrophs microorganism injecting rate	20 mg/L
Oxygen injecting rate	8 mg/L
Water pumping rate	30 L/d
Microorganisms maximum specific growth rate	4.2 per day
Biomass density	0.09 g/cm <sup>3</sup>
Yield coefficient (g cell/g benzene)	1.0 cells/g soil
Half-saturation coefficient	0.77 mg/L
Bulk density of porous medium	1.64 g/cm <sup>3</sup>
Simulation period	12 day

**Table S2.** Initial Geochemical and Microbial Properties of the Soil

Parameter	Value
Soil classification	Silty clay, sand, and clay matrix till
Hydraulic conductivity	In the range of $10^{-7}$ to $10^{-5}$ (m/s)
Moisture content	7.5-32.5% (by volume)
Porosity	30-53.1%
Na	436-548 mg/L
K	16-19.7 mg/L
Ca	562-629 mg/L
Mg	338-407 mg/L
Fe	0.12-1.04 mg/L
Cl	10-79 mg/L
N, NO <sub>2</sub> <sup>-</sup> , NO <sub>3</sub> <sup>-</sup>	31-115 mg/L
Soil organic carbon	1.14%
Dissolved oxygen concentration	<1.0 mg/L to 1.5 mg/L
Initial microbial species	<i>Pseudomonas</i> sp. Strain CFS-215, <i>Geobacter</i> sp., and <i>Rhodococcus</i> sp. Strain 33

**Table S3.** Observed benzene concentrations (mg/L)

Well	Day 13	Day 15	Day 17	Day 19	Day 21	Day 24	Day 26	Day 28	Day 32	Day 34	Day 36	Day 38	Day 40
1		0.032			1.073	0.033		0.174	0.696	0.837	0.738	0.462	0.439
2	0.484	0.564	0.262	0.360	0.672	0.978	0.732	0.699	1.054	1.252	0.682	0.542	0.702
3	0.643	0.734	3.169	2.391	3.408	1.777	2.137	1.858	1.834	1.897	1.077	0.712	0.606
4	0.236				0.245			0.090	0.663	0.827	0.671	0.409	0.415
5	1.347	2.074	1.362	0.888	0.842	0.204	0.601	0.745	0.825	0.974	0.993	0.578	0.685
6	8.131	7.482	7.795	5.530	7.438	7.068	8.696	5.716	4.080	4.887	6.337	3.519	7.340
7	0.279							0.392	0.875	1.370	0.756	0.761	0.566
8	0.296			0.357				0.175	0.851	1.117	0.738	0.594	0.720
9					1.198	0.186		0.300	0.884	1.067	0.724	0.637	0.741
10	1.359	1.218	0.498		1.284	4.698	1.029	1.278	1.384	1.890	1.663	0.546	2.488
11			0.507					0.851	0.508	0.213	0.474	0	0.203
12		0.502	0.808					0.843	0.578	0.288	0.502	0.036	0.352
13													
14													
15													
16									0.485	0.211			0.324

**Table S3.** (continued)

Well	Day 42	Day 44	Day 46	Day 48	Day 52	Day 53	Day 54	Day 57
1	0.304							
2								
3	0.526	1.593	1.429	0.508	0.501	0.733	0.386	0.285
4			0.090	0.400	0.285	0.293	0.292	0.245
5	1.497	1.920	2.070		0.824	0.651	0.581	0.575
6	0.444			0.472	0.265	0.361	0.284	0.291
7	0.385	0.733	1.227	0.686	0.300	0.268	0.241	0.224

8	0.524	0.703	0.366	0.527	0.359	0.316	0.310	0.273
9	1.070	0.918	0.698	1.357	0.831	0.874	0.397	0.633
10	1.366	1.628	0.947	1.590	2.055	1.292	1.162	0.292
11	0.349	0.710	0.163	0.554	0.363	0.280	0.288	0.249
12	0.417	0.741	0.166	0.443	0.322	0.267	0.263	0.285
13			0.079	0.376				
14				0.376	0.231			
15	0.512		0.090	0.398	0.261	0.248	0.237	0.235
16	0.334	0.373	0.453	0.776	0.563	0.507	0.296	0.296

**Table S4.** Error analysis for the biodegradation simulation results

Well number	Observed concentration (mg/L)	Simulated concentration (mg/L)	Absolute Error (mg/L)
3	0.00	0.03	0.03
4	0.00	0.05	0.05
5	0.51	0.25	0.26
6	0.40	0.35	0.05
7	0.80	0.80	0.00
8	0.47	0.16	0.31
9	0.69	0.49	0.20
10	0.53	0.78	0.25
11	1.36	1.70	0.34
12	2.00	2.40	0.40
15	0.41	0.80	0.19
16	0.44	0.20	0.24
Mean absolute error		0.21	
Root mean square error		0.27	
Correlation coefficient		0.93	



**Table S5.** Fifty levels of contamination situation (mg/L)

No.	M 5 $x_1$	M 7 $x_2$	M 8 $x_3$	M 10 $x_4$	M 11 $x_5$	M 12 $x_6$	No.	M 5 $x_1$	M 7 $x_2$	M 8 $x_3$	M 10 $x_4$	M 11 $x_5$	M 12 $x_6$
1	1.68	22.07	5.63	14.10	2.88	1.62	26	6.53	25.45	3.99	2.86	19.79	1.63
2	7.67	4.59	7.35	22.43	22.71	2.05	27	2.34	19.26	20.82	2.65	1.62	1.84
3	8.71	16.50	24.76	4.73	2.57	26.12	28	19.72	5.72	3.73	1.72	3.64	2.33
4	7.53	1.79	12.45	2.26	16.31	20.51	29	9.30	1.93	2.93	4.55	10.89	22.91
5	3.63	4.24	17.74	3.50	2.12	27.21	30	10.25	23.75	24.24	4.31	16.00	2.81
6	14.59	7.04	14.73	10.55	23.17	2.34	31	1.63	1.64	2.38	2.57	5.30	6.46
7	5.27	5.12	6.76	10.95	4.93	20.10	32	5.95	10.31	3.53	3.53	2.17	19.80
8	1.63	16.66	19.79	10.25	16.10	6.22	33	3.71	16.88	14.63	6.29	4.52	3.29
9	3.72	1.89	4.63	13.64	23.44	3.41	34	2.61	14.90	20.27	21.74	2.42	1.92
10	2.17	1.73	6.10	4.88	14.01	5.45	35	6.79	18.64	1.62	12.37	2.39	10.58
11	4.22	1.73	6.59	4.60	2.49	2.37	36	19.48	2.16	27.50	15.27	20.50	3.13
12	10.04	24.28	18.27	10.83	12.57	9.21	37	17.92	2.97	1.64	25.34	2.14	2.30
13	4.99	6.13	25.44	2.06	18.22	2.28	38	4.85	21.68	7.44	26.07	16.37	2.87
14	1.76	21.19	2.29	2.09	2.04	1.72	39	9.17	4.35	3.61	19.46	2.53	23.10
15	2.62	15.15	27.64	3.16	4.30	2.92	40	3.93	19.13	13.06	4.08	23.73	5.44
16	14.17	2.00	15.40	7.13	9.85	13.62	41	14.76	20.89	2.40	19.80	14.06	2.47
17	8.58	27.04	6.77	13.23	15.66	2.10	42	5.04	1.65	21.25	17.10	6.10	17.77
18	15.63	18.38	13.20	6.81	1.68	25.45	43	2.60	4.41	25.02	2.85	1.62	3.82
19	2.22	4.55	7.90	17.96	2.70	2.71	44	4.25	1.73	17.42	1.63	17.27	1.62
20	2.60	15.15	18.04	16.46	25.59	4.85	45	1.63	15.60	22.53	15.09	1.65	23.74
21	2.50	6.36	9.61	4.93	11.25	15.38	46	8.25	1.74	24.30	22.94	1.63	2.68
22	5.41	15.46	5.02	9.83	3.44	22.02	47	3.28	4.63	20.23	4.99	20.29	15.02
23	2.40	5.94	2.23	2.01	12.77	2.09	48	2.42	2.76	1.62	3.05	4.78	6.56
24	4.23	20.07	3.26	16.44	1.83	1.68	49	2.91	1.77	3.79	14.65	8.37	4.08
25	2.21	4.79	2.15	14.61	9.72	2.07	50	19.04	5.04	2.13	1.91	23.67	19.58

**Table S6.** Fifty scenarios of operating conditions

No.	$u_1$ (L/d)	$u_2$ (L/d)	$u_3$ (L/d)	$u_4$ (L/d)	No.	$u_1$ (L/d)	$u_2$ (L/d)	$u_3$ (L/d)	$u_4$ (L/d)
1	37.62	19.42	21.12	13.24	26	10.12	14.16	37.08	10.02
2	15.7	34.46	36.14	30.24	27	11.48	15.66	17.82	27.96
3	34.42	10.84	25.2	11.16	28	13.16	19.38	34.78	21.78
4	14.16	11.48	12.92	38.64	29	20.1	11.08	29.9	18.21
5	32.78	30.18	20.34	12.74	30	23.2	25.1	32.2	12.62
6	27.92	13	11.36	29.72	31	20.5	38.76	10.68	22.76
7	16.64	18.46	32.78	31.4	32	14.6	13.34	19.05	18.76
8	13.08	23.5	30.16	10.08	33	19.8	31.98	26.36	31.88
9	15.6	30.4	13.92	14.92	34	14.92	17.24	38.68	37.88
10	15.6	13.5	29.8	12.78	35	30	27.5	13.6	15.56
11	29.7	26.88	12.6	24.9	36	21.3	12.32	23.36	29.28
12	28.56	29.3	13.26	25.98	37	18.38	16.08	38.64	10.14
13	12.96	15.24	25.28	13.9	38	38	23.6	30.42	32.76
14	35.4	26.68	11.6	13.6	39	17.04	33.84	11.3	14.78
15	26.14	26.98	39.82	20.66	40	19.6	8.42	32.14	38.56
16	27.7	39.78	13.92	11.88	41	29.46	24.22	15.74	23.58
17	28.7	34.9	39.52	12.42	42	16.12	15.06	36.32	33.98
18	34.54	39.78	35.78	22.16	43	16.64	10.66	17.42	34.34
19	14.1	37.58	23.2	13.28	44	33.4	24.4	17.02	17.04
20	24.4	17.66	18.02	38	45	23.32	31.86	28.54	12.72
21	12.12	16.9	20.24	17.32	46	19.06	18.24	11.92	22.44
22	21.34	14.3	35.96	15.78	47	33.96	13.1	32.56	15.09
23	33.42	13.48	23.92	32	48	13.98	13.8	39.72	39.6
24	22.5	36.2	13	22.64	49	20.06	34.96	16.54	38.44
25	20.2	27.94	34.72	25.84	50	39.8	26.42	38.82	21.52

**Table S7.** Input and output variables for SI emulator and FCI simulator

SI Emulator	Input (I) or Output (O)	Symbol	FCI Simulator	Input (I) or Output (O)	Symbol
Highest contaminant concentration anywhere in the mesh	I	$\bar{b}^{\text{MAX}}$	Highest contaminant concentration anywhere in the mesh	I	$\bar{b}^{\text{MAX}}$
Percentage of benzene mass removal	I	$\eta$	Percentage of benzene mass removal	I	$\eta$
Injecting rate of well I	O	$u_1$	Injecting rate of well I	I	$u_1$
Injecting rate of well II	O	$u_2$	Injecting rate of well II	I	$u_2$
Pumping rate of well III	O	$u_3$	Pumping rate of well III	I	$u_3$
Pumping rate of well IV	O	$u_4$	Pumping rate of well IV	I	$u_4$
			Highest contaminant concentration anywhere in the mesh	O	$\bar{b}^{\text{MAX}}$
			Percentage of benzene mass removal	O	$\eta$

**Table S8.** The result of the collinearity test

Variable	VIF	Tolerance
u1	1.068	0.936
u2	1.128	0.886
u3	1.043	0.959
u4	1.067	0.937