

Physics-based Pollutant Source Identification in Stormwater Systems

Andrew Chio*, Russell Bent[†], Andrey Y. Lokhov[†], Jian Peng[‡] and Nalini Venkatasubramanian*

*Dept. of Computer Science, University of California, Irvine, {achio,nalini}@uci.edu

[†]Theoretical Division, Los Alamos National Laboratory, {rbent,lokhov}@lanl.gov

[‡]Orange County Public Works, jian.peng@ocpw.ocgov.com

Abstract—Stormwater networks are critical utility infrastructures designed to drain rainwater and nuisance flows, such as excess irrigation and groundwater seepage from urban communities. During this process, they can transport pollutants (e.g., pesticides, oils, and greases) to receiving waters such as rivers, bays and oceans. A recurring problem faced by these systems are *dry weather flows* (DWFs), where illicit discharges are introduced and propagated in the network during periods with no rain. Current techniques for monitoring DWFs consist of manual inspections and grab samples, which are costly and inefficient. However, with advances in sensing and communication, the Internet-of-Things (IoT) has enabled new opportunities for enhanced decision support and control. This paper proposes a quick and efficient physics-based backwards inference model to identify potential sources of pollutant discharges in DWFs, given time-series IoT observations and knowledge embedded in domain-expert simulations. Our approach leverages the underlying physics that drives flow propagation in stormwater systems, and optimizes multiple least-squares regressions to find potential DWF sources and their associated flows. We evaluate our backwards inference model on six real-world stormwater networks provided by domain experts, and show its efficacy in reconstructing anomalies.

I. INTRODUCTION

Stormwater networks, also known as municipal separate storm sewer systems (MS4s), are a utility infrastructure consisting of thousands of catch basins, outfalls, and channels. The system transports rainwater and nuisance flows (e.g., excess irrigation, groundwater seepage) from urban cities to receiving waters such as rivers, bays, and oceans. However, these networks can also carry pollutants like pesticides, oils, and sewage, which may be introduced by illegal dumping, incidental discharge or illicit non-MS4 pipeline connections. This can lead to water quality impairments downstream, such as increased harmful bacteria, algal blooms, and fish kills.

To address this, regulations like the 1987 amendment of the US Clean Water Act [1] require permits for urban MS4 discharges, and generally prohibit non-MS4 connections and their associated pollutants. However, enforcement against violators and remediation is challenging due to the difficulty of tracing discharge origins, which could be transient in nature and occur over large catchment areas. In this context, source identification and mitigation efforts like [2], [3] focus on *dry weather flows* (DWFs), which occur during dry weather/low flow conditions with no rain. These discharges are easier to identify, and managing them also improves water quality during wet weather.

Several aspects of DWF discharges make source identification difficult. One key challenge is the transient and sudden

nature of DWFs, where pollutants can appear abruptly, cause significant harm, and dissipate before source tracking and clean up measures can be initiated. Another challenge comes from the large, geo-distributed scale and complex nature of MS4s. For instance, catch basins in typical cities lie every 100-150m, and connect to increasingly larger underground storm pipes before daylighting. Thus, DWFs from small urban drainage areas can originate from any of thousands of potential points, making it hard to quickly identify pollutant sources and take corrective measures.

Current state-of-the-art methods for monitoring DWFs are inadequate for enforcing stormwater permits [2], [4], [5]. This typically involves visual inspections, citizen reports, and water quality sampling, where domain experts employ test kits and lab analysis. For example, Orange County Public Works conducts 5 site visits each for 30 selected stormwater outfalls each summer, where approximately an hour was spent for observation, testing, and sampling at each site. These practices are costly, time-consuming, and ineffective: capturing an ongoing illegal discharge event is extremely unlikely; and lab results, which require 3-5 weeks to process, are impractical to act upon. However, with the rise of water quality and flow sensing technologies, the Internet of Things (IoT) has shown great promise in enabling low cost, real-time and continuous monitoring, and quick analysis capabilities, which can support effective decision-making and control.

This paper focuses on the problem of *DWF source identification*, where pollution sources of DWFs within MS4s are inferred using observations from pre-deployed IoT sensors. Through this analysis, stormwater agencies can start to reliably detect pollution incidents and enforce stormwater permit regulations. To understand the evolution of DWFs, domain expert simulations are typically run and cached to estimate potential network states. This can be computationally expensive and memory intensive, making it impractical for real-time analysis on large, geo-distributed networks where quick decision support and control is essential. We address these issues by leveraging the physics of flow propagation to gain insight for a quick and efficient backwards inference model. The key contributions of this paper are as follows:

- A review of related literature for pollution source identification and the role of physics in simulations (§II)
- A physics-based backward inference model for potential DWF origins, given IoT sensor observations. *A salient feature of this contribution is a physics approximation that balances model fidelity with the computational requirements*

of backwards inference. (§III and §IV)

- An evaluation of our model with six real-world stormwater networks and discussion on future directions (§V, §VI)

II. RELATED WORKS

In this section, we examine related works for techniques of source identification in networked systems, and the role that physics-based models play for inference.

Source Identification in Networked Systems. The problem of anomaly source identification has been studied in many different networked systems domains, from smart grids [6], to healthcare [7]. In these settings, research efforts have explored the modeling of the underlying systems and their operations to study the different types of impacts that an anomaly could have in the network. In this paper, we focus on source identification as it applies to the water domain. Techniques to identify the source of an anomaly are rich and varied in the water domain. Traditional methods rely on physically observing, sampling, and post hoc lab analysis to identify a water quality anomaly [2], [8], [9]. However, low probability of encountering an observable anomaly in the field (less than 0.01 percent) and delays in obtaining results and uncertainties in pinpointing potential upstream origins, makes source identification oftentimes infeasible.

In general, given the dynamic and unpredictable nature of the anomalies and complexity OF the MS4, solutions for source identification are usually not unique. A common approach is the Bayesian approach, where unknown variables (e.g., the pollutant origin) are treated as random variables, and iteratively updated using data from sensors. The Bayesian approach was first applied to the drinking water setting in [10] to estimate a contaminant’s release history; they considered a formulation that only used the distance from the source. This was extended in [11] to address mixed uncertainty in models; they cast the source identification problem as a minimax problem. However, Bayesian methods can be expensive to run, since they update probabilities in a continuous manner using distributions of variables. Later works, such as [12], [13] have looked to address this issue by considering a more sparse input, and locally accurate approximations using Markov Chain Monte Carlo methods. Other than Bayesian methods, optimization-based approaches have also been proposed, which utilize greedy heuristics [14] and evolutionary algorithms [15], [16]. This typically provides faster computation times, at the cost of yielding sub-optimal solutions. Machine learning [17] and other recent advances in deep learning [18], [19] also show promise, but generally require large amounts of data, combined with heavy manual tuning and computation to produce accurate solutions.

Role of physics-based models. The nature and behavior of flows and pollutants in the interconnected MS4 systems are governed by fundamental laws of physics, such as the conservation of mass and energy. This provides a basis on which domain expert simulations are built, and run to gain insight into potential future or past events. Examples include the start of an anomaly, and a rapid changes in network

states. Source identification approaches can be categorized by how the underlying physics and simulation models are used. Many approaches adopt a “black box” methodology, where physics-based simulations are run on a predefined set of inputs and cached [20]–[23]. This reduces the source identification problem to a search for an appropriate sample input. However, the accuracy and reliability of these methods depend on the quality of their inputs, and can require massive compute and memory resources. Other works involve physics in a “white box” manner, such as in [24]–[26]. While this exposes the complexity of the embedded physics and necessitates deep knowledge of the domain, the underlying computational model can be exploited to study the effect of specific inputs. Our work is set in the context of the latter case, and aims to construct a backwards inference model for stormwater networks *by deriving a close approximation of the physics to directly interact with governing equations of stormwater flow dynamics.*

III. PROBLEM FORMULATION

We first present models of the stormwater infrastructure network, and the physics of flow propagation. Then, we define dry weather flow (DWF) anomalies and sensors.

Stormwater Infrastructure Network. Let the directed acyclic graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ denote the geo-distributed stormwater network, with junction nodes \mathcal{V} which represent potential DWF insertion points (e.g., catch basins), and conduit edges $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$ which represent pipes that propagate flow between nodes via gravity. Each junction $v_j \in \mathcal{V}$ is characterized by its physical location (x_j, y_j) and invert elevation z_j . Each conduit $e_{ij} \in \mathcal{E}$ is characterized by its length L_{ij} , frictional roughness f_{ij} , and cross-sectional shape S_{ij} , and its slope is derived using the elevations of its ends. These attributes are used to simulate the physics of network flow, which carries the pollutants of interest.

Physics of Flow Propagation. The propagation of flow in stormwater networks is governed by the conservation of mass (Eq.1a) and momentum (Eq.1b). These equations relate many physical quantities: distance x , time t , flow area A , flow rate Q , hydraulic head H , friction slope S_f and gravity g .

$$\frac{\partial A}{\partial t} + \frac{\partial Q}{\partial x} = 0 \quad (1a)$$

$$\frac{\partial Q}{\partial t} + \frac{\partial(Q^2/A)}{\partial x} + gA \frac{\partial H}{\partial x} + gAS_f = 0 \quad (1b)$$

These equations are often solved using Euler’s method and the United States Environmental Protection Agency Storm Water Management Model (EPA SWMM) [27] is one the most widely used implementations. This open-source stormwater simulator models and solve for several complex and time-varying interactions between runoff, flow routing, DWFs, backwater effects, losses and more. This is accomplished through a *dynamic wave analysis* method which solves Eq. 1a and 1b by expressing them as functions of Q and H . Then, finite approximations are made and the implicit backwards Euler method is applied for each junction

and conduit. This is a discrete process that iteratively finds values for Q and H through the following update functions:

$$Q^{t+\Delta t} = \frac{Q^t + \Delta Q_{iner} + \Delta Q_{pres}}{1 + \Delta Q_{fric}} \quad (2a)$$

$$\Delta Q_{iner} = 2\bar{U}(\bar{A}^{t+\Delta t} - \bar{A}^t) + \bar{U}^2 \frac{(A_{dn} - A_{up})}{L} \Delta t \quad (2b)$$

$$\Delta Q_{pres} = -g\bar{A} \frac{(H_{dn} - H_{up})}{L} \Delta t \quad (2c)$$

$$\Delta Q_{fric} = g\eta^2 \frac{|\bar{A}|\Delta t}{R^{4/3}} \quad (2d)$$

$$H^{t+\Delta t} = H^t + \frac{\Delta t/2(\sum Q^t + \sum Q^{t+\Delta t})}{(A_{SN} + \sum A_{SL})^{t+\Delta t}} \quad (2e)$$

The flow rate $Q^{t+\Delta t}$ in Eqn. (2a) is updated through changes in inertia ΔQ_{iner} , pressure ΔQ_{pres} , and pipe friction ΔQ_{fric} , as shown in Eqn. (2b), (2c), and (2d), respectively. These quantities depend upon the hydraulic head at the upstream and downstream ends (denoted as H_{up} and H_{dn}), and the flow area at the upstream and downstream ends (denoted as A_{up} and A_{dn}). The updated flow rate values are then used to compute the hydraulic head for the next timestep in Eqn. (2e). The fidelity of this model is further enhanced by other details described in [27] and Appx. A. However, this makes EPA SWMM difficult to use for inference as it introduces non-differential equations. Thus, we apply several approximations to achieve differentiability.

Approximations for Differentiability. Here, we describe key approximations to the dynamic wave analysis method in EPA SWMM [27], so that a computational graph can be derived with automatic differentiation, and used for inference. Additional computational details can be found in Appx. A.

First, we remove the boundary conditions to identify whether a node is considered “dry”, i.e., no flow at node. This is done to allow the computational graph to be derived without introducing a dependence on the value of a variable. While this can create discrepancies where dry nodes in the simulation have small flows, our experiments in §V show its negligible impact for inference applications. We also remove limitations on critical and surcharged flows in the network, which occur when pipes are (intuitively) full; this does not realistically occur for our driving use case with DWF anomalies, as we consider periods without rainfall. By extension, support for backflow propagation in pipes is removed, since this can only occur after a pipe is surcharged. Next, we increase the number of iterations used to converge to a stable solution to a constant, instead of terminating early. Lastly, discrete functions, e.g., weight factors measuring the super-criticality of flows, are approximated with continuous, differentiable equations. Through these modifications, the simulation becomes differentiable, enabling fast nonlinear solvers to be used for backwards inference. Our implementation is found on GitHub [28].

DWF Anomalies. We define a DWF anomaly $\alpha_k \in \mathcal{A}$ as an illicit discharge event that introduces pollutants into the stormwater network over time. We characterize each anomaly by its origin node v_k^* representing a source location, and

a DWF inflow curve $Q_{v_k^*}^{dwf}(t)$, representing the amount of flow introduced at v_k^* at time t . This inflow is bound by $[Q^{min}, Q^{max}]$, representing the minimum and maximum inflows at v_k^* . We assume that $Q_{v_k^*}^{dwf}$ has a limited-time injection profile, in order to mimic a potential illegal dumping. This makes anomalies transient in nature: they have a limited time in the network over which detection is possible, as shown in Fig. 1. We note that the problem of *anomaly detection*, where anomalies are identified by a rapid change in water quality parameters (e.g., pH, temperature, electric conductivity and turbidity) are fundamentally related to, and propagated by, flow in the network. However, this is out of the scope of this paper, and left as a future extension.

Sensors. A flow sensor $s_l \in \mathcal{S}$ measures the flow rate at an instrumented node with a periodicity of λ_l seconds. We denote the flow rate data captured by s_l by the time series $\{Q_{s_l}^{obs}(t)\}_{t \in T}$. Fig. 1 shows the observations made by a sensor downstream, from which our goal is to reconstruct the upstream flows, which we detail next.

IV. METHODS

In this section, we present the general overview of our DWF backwards inference model, which relies on optimizing a least-squares regression over a set of candidate nodes.

Problem Statement and General Approach. Consider a DWF anomaly α_k that introduces an unknown amount of flow $Q_{v_k^*}^{dwf}$ at node v_k^* . Then, suppose that $\mathcal{S}^* \subseteq \mathcal{S}$ is the set of sensors deployed in the network that observe α_k . Let $\{Q_{s_l}^{obs}\}_{s_l \in \mathcal{S}^*}$ denote the flow rate observations made at times $t \in T$. The *source identification* problem then infers the flow $Q_{v^*}^{inf}$ to be introduced at node v^* that would most likely produce $\{Q_{s_l}^{obs}\}_{s_l \in \mathcal{S}^*}$, for a set of candidate nodes \mathcal{V}^* . For this, we leverage the differentiable physics model in §III, which produces a set of simulated edge flows $\{Q_{s_l}^{simu}\}_{s_l \in \mathcal{S}^*}$ from an anomaly at v^* . Fig. 1 illustrates this problem: given sensor observations, construct the inferred flows (orange) that would best match the ground truth flows (blue).

Our approach begins by identifying a set of potential sources \mathcal{V}^* that could reasonably produce the set of observations $\{Q_{s_l}^{obs}\}_{s_l \in \mathcal{S}^*}$. For each identified source $v^* \in \mathcal{V}^*$, we formulate a least squares regression whose solution provides an associated DWF flow $Q_{v^*}^{inf}$ that would need to be introduced at v^* to generate the observations $\{Q_{s_l}^{obs}\}_{s_l \in \mathcal{S}^*}$. We note that the solution to the source identification problem is oftentimes not unique, as several equally-likely origins and corresponding flows could produce the same set of observations. In this paper, we focus on identifying all such potential sources (and their associated flows), which can then be provided to a domain expert for further analysis. We assume that sensors are already deployed in the network, and that flow is a approximate indicator of the pollutants and contaminants of interest.

Finding potential source nodes. Our backwards inference first constrains the search space of potential source nodes. Since gravity is the main force propagating flow between nodes in a stormwater network, we can limit the set \mathcal{V}^* to

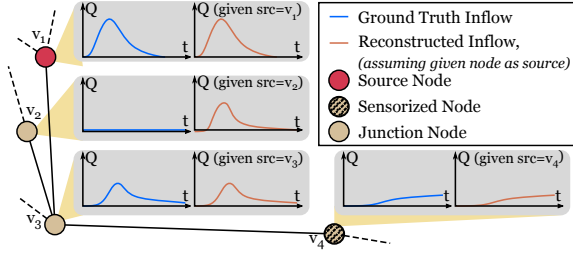


Fig. 1. An example of flow propagation in a sample network

those that lie upstream of all sensors observing the anomalous flow. Moreover, we note that the lack of observations at a sensor can also provide insight into eliminating impossible origin nodes. However, due to the transient nature of the DWF anomalies, distinguishing between the absence of an anomaly, and one that cannot be detected (e.g., too low flow) without observations is impossible. Therefore, we use a distance threshold τ_d to ensure that the nodes removed are local to a sensor. Note that the selection of τ_d is closely tied to the expected transient behavior of anomalies in the network. Intuitively, if τ_d is large, then our approach caters to DWF anomalies that occur over longer timescales and/or have larger impacts in the network. On the other hand, when τ_d is small, our approach increasingly considers very transient anomalies with small DWF injection profiles. For both cases, however, a number of potential source nodes could be falsely eliminated. Future work will identify suitable values for τ_d by leveraging domain expert knowledge and feedback, with insights extracted from historical data.

We detail the general process of finding potential source nodes in the first part of Alg. 1. We start by initializing \mathcal{V}^* to the set of all nodes in the graph \mathcal{G} (line 1). We then iterate through each deployed sensor $s_l \in \mathcal{S}$, and obtain the set of nodes that lie upstream of s_l 's location within distance τ_d (lines 2-4). The set of potential nodes is then constrained based on if s_l observed the anomaly (line 5-6), or not (line 7-8). We describe the rest of the algorithm later.

Formulating the Least Squares Regression. For each candidate origin node v^* , we construct a least squares regression that attempts to match the set of ground truth flows in the network. We simulate the values $\{Q_{s_l}^{simu}\}_{s_l \in \mathcal{S}}$, which were computed under the assumption that v^* is the origin. Then, the least squares optimization problem is:

$$\underset{Q_{v^*}^{dwf}}{\operatorname{argmin}} \sum_{s_l \in \mathcal{S}^*} \sum_{t \in T} \left(Q_{s_l}^{obs}(t) - Q_{s_l}^{simu}(t; Q_{v^*}^{dwf}) \right)^2 \quad (3a)$$

$$\text{s.t.} \quad Q^{min} \leq Q_{v^*}^{dwf}(t) \leq Q^{max} \quad \forall t \in T \quad (3b)$$

$$\text{Eqn. (2a) to (2e)} \quad \forall t \in T \quad (3c)$$

The regression objective in Eqn. (3a) minimizes the difference between observed flow values $Q_{s_l}^{obs}(t)$ and the theoretical physics-based flow values $Q_{s_l}^{simu}(t; Q_{v^*}^{dwf})$, across the sensors observing the flow $s_l \in \mathcal{S}^*$ for each time step $t \in T$. The constraint in (3b) limits the search space for the injected flow, and the constraints of (3c) dictate how $Q_{s_l}^{simu}(\cdot)$ is computed based on the dynamics of the system.

Algorithm 1: Source Inference

Input: Graph \mathcal{G} , Sensors \mathcal{S} , Observations $\{Q_{s_l}^{obs}\}_{s_l \in \mathcal{S}}$,
float τ_q , float τ_d , float τ_o
Output: Potential anomalies \mathcal{A}^*

```

// Compute set of potential origin nodes
1  $\mathcal{V}^* \leftarrow \emptyset$ 
2 for  $s_l \leftarrow \mathcal{S}$  do
3    $v \leftarrow \text{GetDeployedNode}(s_l)$ 
4    $\mathcal{V}^{up} \leftarrow \text{GetUpstreamNodes}(v, \mathcal{G}, \tau_d)$ 
5   if  $\exists t : Q_{s_l}^{obs}(t) \geq \tau_q$  then
6     //  $s_l$  observed the anomaly at some time
7      $\mathcal{V}^* \leftarrow \mathcal{V}^* \cup \mathcal{V}^{up}$ 
8   else
9     //  $s_l$  did not observe the anomaly
10     $\mathcal{V}^* \leftarrow \mathcal{V}^* - \mathcal{V}^{up}$ 
// Find flow curves for potential sources
9  $M \leftarrow \text{map}()$ 
10 for  $v^* \leftarrow \mathcal{V}^*$  do
11    $objval, Qinf \leftarrow \text{solve Eqn. (3), assuming source } v^*$ 
12    $M[src] \leftarrow (objval, Qinf)$ 
13  $minobjval \leftarrow \min_{objval, \cdot \in M.values()} \{objval\}$ 
14  $\mathcal{A}^* \leftarrow \{(src, Qinf) : \forall (src, (objval, Qinf)) \in M : |minobjval - objval| \leq \tau_o\}$ 
15 return  $\mathcal{A}^*$ 

```

The second half of Alg. 1 explains how our least squares regression is leveraged for backwards inference. After obtaining the set of potential source nodes \mathcal{V}^* , we construct the associated least squares regression for each node $v^* \in \mathcal{V}^*$ (lines 9-12). We note that this step can be parallelized for improved computation time for inference. The last step looks to pick sources nodes and accompanying DWF flows that could all potentially realistically occur - we return all such solutions that are within a threshold τ_o of the minimum objective value found (lines 13-15).

V. EXPERIMENTS AND RESULTS

In this section, we evaluate our physics-informed backwards inference model on six real-world stormwater networks. Our experiments examine the impacts of the approximations, and the quality of the resulting inferences.

Experimental Setup. Our evaluation leverages six real-world stormwater networks provided by Orange County Public Works. These networks are defined using EPA SWMM [27], and consist of: three small networks of ~ 350 to 700 nodes and edges over a ~ 100 km² area (Fig. 2(a), 2(b), 2(c)); two medium networks of ~ 1000 nodes and edges over a ~ 200 km² area (Fig. 2(d), 2(e)); and one large network of ~ 1500 nodes and edges over a ~ 400 km² area (Fig. 2(f)). To fit with our implementation of the backwards inference model, we modify all network conduits to consist of rectangular pipes; we leave the extension of our model to different cross-sections as future work.

In these networks, we assume that homogeneous flow sensors that generate observations with periodicity $\lambda=30s$ are deployed. Our experimental results are obtained wrt. a set of 100 anomalies that were constructed randomly for each network: each anomaly was assigned a randomly chosen origin node, with an inflow curve that has a peak

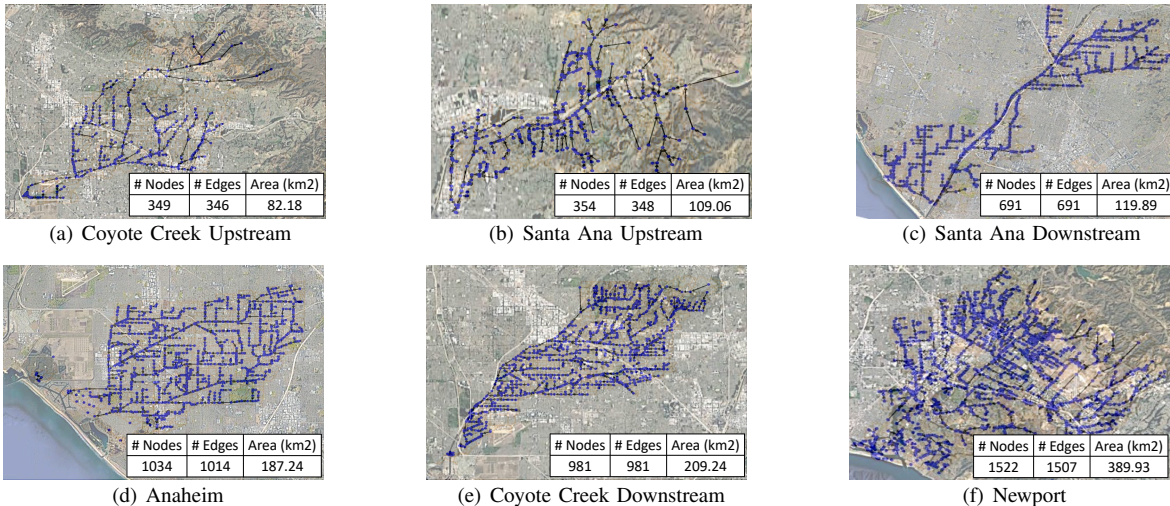


Fig. 2. EPA SWMM Networks used for Evaluation

flow magnitude of $[0.25 \pm 0.2]$ cfs, and start and end times chosen randomly between 0 and 2 hours. Our experiments were conducted on an M1 MacBook Pro with 16 GB of memory and 10 CPU cores. Our backwards inference model is implemented in Julia, using the JuMP [29] interface to Ipopt [30] and the MA57 linear solver [31] for optimization. We publish our backwards inference model on GitHub [28].

Comparison Baseline. We compare our backwards inference model with a standard “black-box” approach [20]–[23] to source identification, which caches edge flows across a large set of predefined anomalies and “infers” an anomaly by searching for the best match to a given set of sensor observations. We simulated and cached 10 anomalies uniformly across all junctions of each network, for a total of ~ 3000 – 15000 cached anomalies for each network, depending on its size. Corresponding peak flow and duration parameters were chosen in the same manner as the set of anomalies used for evaluation.

A. Impact of Approximations for Differentiability

Our first experiment examines the impact of the approximations on the accuracy of modeling and solving for flow dynamics. To this end, we simulate edge flows produced by EPA SWMM for each anomaly created for evaluation. This is then compared with our differentiable version, and the mean square error (MSE) between edge flows is reported, as seen in Table I. In order to avoid artificially improving the MSE, we exclude edges that were not impacted by the anomaly using a minimum flow threshold of 0.0001 cfs.

TABLE I
IMPACT OF APPROXIMATIONS

Network	Avg MSE
Anaheim	$2e-3 \pm 2e-2$
Coyote Creek Downstream	$9e-6 \pm 1e-5$
Coyote Creek Upstream	$3e-4 \pm 2e-3$
Newport	$3e-6 \pm 4e-6$
Santa Ana Downstream	$1e-2 \pm 8e-4$
Santa Ana Upstream	$1e-5 \pm 1e-5$

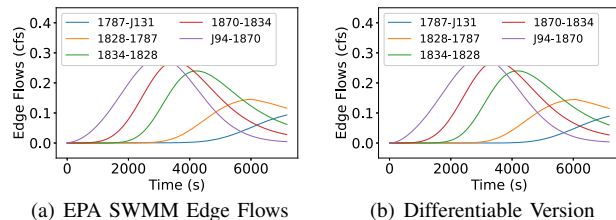


Fig. 3. Comparison of edge flows from a typical anomaly

This shows that the approximations made were negligible: the differentiable variant closely models and solves for flow dynamics as EPA SWMM does. This is illustrated in Fig.3, which plots a typical anomaly simulated using EPA SWMM (Fig.3(a)), and our differentiable version (Fig.3(b)). Note that inference results for this anomaly did not change with 75% or 100% instrumentation, and so these lines are omitted. Thus, this key contribution of our work is suitable for inference.

B. Evaluation of the Backwards Inference Model

We next examine the accuracy and time taken to reproduce an anomaly’s flow curve using sensor observations. We consider varying levels of instrumentation in the network, ranging from 10% to 100% of the nodes having sensors. Fig. 4 reports the MSE between each of the 100 evaluated anomaly’s ground truth flows, and our model’s corresponding inferred flows for the true origin of the anomaly. Our inference model was able to reconstruct the anomaly inflow with little error: the average MSEs across the small, medium, and large networks was 0.02, 0.018, and 0.023, when 10% of the network was instrumented, and all decrease to ~ 0 as the level of instrumentation increased. This is due to the redundancy of detection, which allows our model to optimize the inferred inflow to be consistent with all captured observations. In comparison, the baseline standard (depicted with colored dashed lines for each network) follows the same trend, but generally performs worse than our inference model, and only marginally improves as the level of instrumentation increases. We show these trends using the thicker black lines in Fig. 4, which plots the average MSEs over all

networks. The standard deviation averaged across networks is represented by the tan region (our inference model) and light green region (baseline), which decreases as more sensors were instrumented. For our inference model, this ranges from 0.016–0.032 when 10% of the network was instrumented, to 0.0–0.01 when the network is fully instrumented, and is always smaller than that of the baseline standard. To visualize these MSEs, Fig. 5 plots the ground truth and inferred inflows for different levels of instrumentation for the anomaly plotted in Fig 3, as well as the best baseline result with full instrumentation.

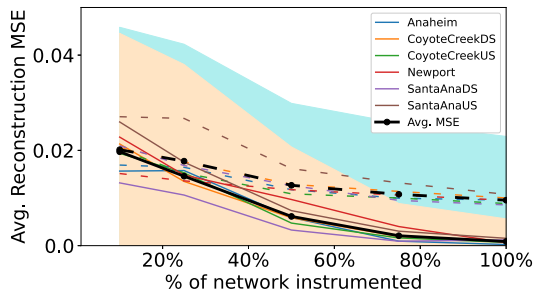


Fig. 4. MSE in Backwards Inference

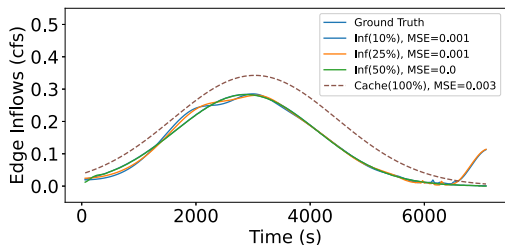


Fig. 5. Comparison of Inferred Flows from a typical anomaly

We also report the average inference time for our model in Fig. 6. Our results show that the level of instrumentation is directly proportional with the inference time. On average, with 10% instrumentation, it took 256 ± 129 s, 248 ± 8 s, and 258 s to produce the inferred flow for the small, medium, and large networks, respectively, and increases to 819 ± 171 s, 1170 ± 64 s, and 1291 s with full network instrumentation. While the time to search cached values is significantly less than the computation needed for inference, we note the tradeoff in the accuracy of results, as well as the storage and offline time needed. In particular, producing this cache across the six networks took ~ 3 days, and used ~ 14.3 GB of memory, which can become prohibitive to run as the network grows. Thus, we show that our model is able to reproduce anomaly flows both accurately and quickly, which is essential in supporting real-time control for managing anomalies.

C. Degeneracy of Inference Results

Lastly, we report the degree of degeneracy in the inference results. As mentioned in §IV, the source identification problem is difficult due to the indistinguishability between potential upstream anomaly origins. Thus, it is critical to reduce the number of potential sources reported, to enable practical decision support for domain experts and practitioners.

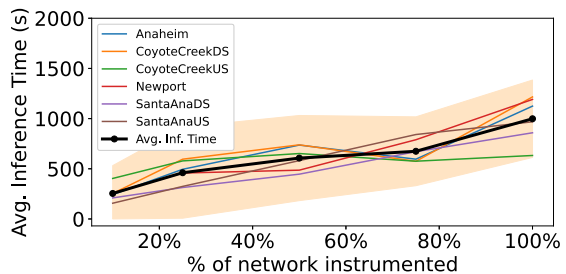


Fig. 6. Time Taken in Backwards Inference

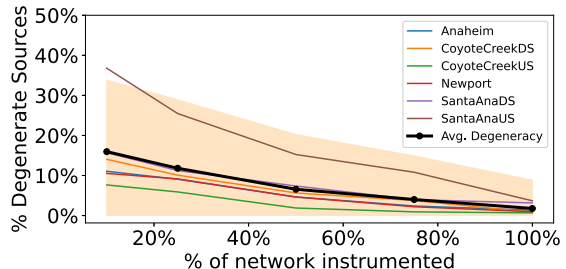


Fig. 7. Degeneracy in Backwards Inference

In Fig. 7, we report the number of other equally-likely potential sources (and corresponding inflows) that were found by our inference model, apart from the true anomaly. This represents nodes at which a specific inflow could result in downstream edge flows consistent with sensor observations. As the number of observations increases in the network, the degree of uncertainty on the source of the anomaly decreases. We note that other methods of reducing the size of the degenerate set include introducing prior probabilities on the potential source nodes, which we leave as future work.

VI. CONCLUSIONS AND FUTURE WORK

In this paper, we presented a physics-informed backwards inference model for stormwater infrastructure systems. Our approach identifies the physical properties of an anomaly in a two step process. First, we use a constraint-based method to eliminate unlikely anomaly source nodes. Then, a least squares regression problem was formulated, whose solution would infer a potential injection profile. We then enabled quick and efficient optimization of this regression using fast nonlinear solvers, combined with a unique set of approximations made on the modeling and solving of stormwater flow dynamics. Our experiments leveraged six real-world stormwater networks, and showed the negligible impacts of the approximations, as well as the efficacy of the backwards inference, even in the case of partial network observability. Future directions include the integration of this model into an end-to-end system for source identification in stormwater networks. In this respect, it is expected that network topology and domain expert feedback will further enrich the observed sensor data. We aim to more closely analyze the observability of the internals of the system to better understand the sensitivity of the different components of the system. Additionally, we will examine the scalability and robustness of the model to incorporate better approximations of the physics, and applications towards other types of networks and anomalies.

DETAILS ON APPROXIMATIONS FOR DIFFERENTIABILITY

Here, we describe the details concerning the approximations made to enable differentiability in our approach. We note that Eqn. (2) expresses the main discretization of the fundamental physics in Eqn. (1a) and (1b). However, in practice, EPA SWMM [27] employs several additional error correction computations which make it difficult to infer the computational graph.

Our first approximation removed boundary conditions for negligible flows. We note that the EPA SWMM implementation follows a scheme closer to that of Eqn. (4), where \bar{A} is the average flow area at the node, and ϵ_A is a minimum flow area threshold. Thus, our implementation only considers Eqn. (2a).

$$Q^{t+\Delta t} = \begin{cases} \text{Eqn. (2a)} & \text{if } \bar{A} \geq \epsilon_A \\ Q^t & \text{else} \end{cases} \quad (4)$$

Critical and surcharged flows only occur when flow in a conduit reaches or exceeds its maximum “capacity”. Several conditions are applied in EPA SWMM which dictate whether a normalized flow should replace $Q^{t+\Delta t}$. These conditions check that: (i) $Q^{t+\Delta t} > 0$; (ii) the conduit is not already in a critical state, nor is it under an edge case with one critical and one dry end; (iii) the upstream flow velocity exceeds a critical velocity. Eqn. (5) shows the normalized flow used under these conditions, where n is the Manning roughness coefficient, and R is the cross-sectional flow radius.

$$Q^{t+\Delta t} = \min \left\{ Q^{t+\Delta t}, Q_{norm} \right\}, \text{ where} \quad (5)$$

$$Q_{norm} = \frac{1.49}{n} A_{up} R_{up}^{2/3} \sqrt{L^2 - (H_{up} - H_{dn})^2}$$

Next by extension, since we assume that critical and surcharged flows do not occur in the network, we additionally discard a postprocessing step in EPA SWMM that modifies the sign of $Q^{t+\Delta t}$, as shown in Eqn. (6).

$$Q^{t+\Delta t} = \begin{cases} Q^{t+\Delta t} & \text{if } Q^t \cdot Q^{t+\Delta t} > 0 \\ 0.001 * \text{sign}(Q^{t+\Delta t}) & \text{else} \end{cases} \quad (6)$$

Lastly, EPA SWMM leverages piecewise-defined weight factors that measure the closeness to criticality in conduits. These weight factors are used in modifying the average flow area and radius to be more numerically stable, i.e., $\bar{A}' = A_{up} + \sigma(\bar{A} - A_{up})$ and $\bar{R}' = R_{up} + \sigma(\bar{R} - R_{up})$. Our approach replaces the computation of σ with Eqn. (7).

$$\sigma = \left(1 + \exp \left(10 * \left(|\bar{U}| / \sqrt{g\bar{A}\bar{W}} - 0.75 \right) \right) \right)^{-1} \quad (7)$$

ACKNOWLEDGEMENTS

This work is supported by the UC National Laboratory Fees Research Program Grant No. L22GF4561, and NSF Grant No. 1952247. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

- [1] C. Copeland, “Clean water act: a summary of the law.” Congressional research service, Library of Congress Washington, DC, 1999.
- [2] B. Bernstein *et al.*, “Assessing urban runoff program progress through a dry weather hybrid reconnaissance monitoring design,” *Environ. Monit. Assess.*, vol. 157, 2009.
- [3] K. Halbach *et al.*, “Small streams—large concentrations? pesticide monitoring in small agricultural streams in germany during dry weather and rainfall,” *Water Research*, vol. 203, 2021.
- [4] A. Barbosa *et al.*, “Key issues for sustainable urban stormwater management,” *Water research*, vol. 46, no. 20, 2012.
- [5] M. K. Leecaster *et al.*, “Assessment of efficient sampling designs for urban stormwater monitoring,” *Water research*, vol. 36, no. 6, 2002.
- [6] H. Jiang *et al.*, “Fault detection, identification, and location in smart grid based on data-driven computational methods,” *IEEE Trans. Smart Grid*, vol. 5, no. 6, 2014.
- [7] H. Rossman *et al.*, “A framework for identifying regional outbreak and spread of covid-19 from one-minute population-wide surveys,” *Nature Medicine*, vol. 26, no. 5, 2020.
- [8] H. Lee *et al.*, “Design of stormwater monitoring programs,” *Water research*, vol. 41, no. 18, 2007.
- [9] D. Li *et al.*, “Municipal separate storm sewer system (ms4) dry weather flows and potential flow sources as assessed by conventional and advanced bacterial analyses,” *Environmental Pollution*, 2023.
- [10] M. F. Snodgrass *et al.*, “A geostatistical approach to contaminant source identification,” *Water Resour. Res.*, vol. 33, no. 4, 1997.
- [11] A. Y. Sun, “A robust geostatistical approach to contaminant source identification,” *Water Resour. Res.*, vol. 43, no. 2, 2007.
- [12] L. Zeng *et al.*, “A sparse grid based bayesian method for contaminant source identification,” *Advances in Water Resources*, vol. 37, 2012.
- [13] J. Zhang *et al.*, “Efficient bayesian experimental design for contaminant source identification,” *Water Resour. Res.*, vol. 51, no. 1, 2015.
- [14] B. K. Banik *et al.*, “Greedy algorithms for sensor location in sewer systems,” *Water*, vol. 9, no. 11, 2017.
- [15] M. M. Aral *et al.*, “Identification of contaminant source location and release history in aquifers,” *J. Hydrol. Eng.*, vol. 6, no. 3, 2001.
- [16] K. Han *et al.*, “Application of a genetic algorithm to groundwater pollution source identification,” *J. Hydrol.*, vol. 589, 2020.
- [17] L. Grbčić *et al.*, “Water supply network pollution source identification by random forest algorithm,” *J. Hydroinform.*, vol. 22, no. 6, 2020.
- [18] S. Mo *et al.*, “Deep autoregressive neural networks for high-dimensional inverse problems in groundwater contaminant source identification,” *Water Resources Research*, vol. 55, no. 5, 2019.
- [19] A. Solanki *et al.*, “Predictive analysis of water quality parameters using deep learning,” *IJCA*, vol. 125, no. 9, 2015.
- [20] J. Guan *et al.*, “Identification of contaminant sources in water distribution systems using simulation—optimization method: case study,” *J. Water Resour. Plan. Manag.*, vol. 132, no. 4, 2006.
- [21] P. Vankayala *et al.*, “Contaminant source identification in water distribution networks under conditions of demand uncertainty,” *Environmental Forensics*, vol. 10, no. 3, 2009.
- [22] Q. Han *et al.*, “Aqueis: Middleware support for event identification in community water infrastructures,” in *ACM Middleware*, 2019.
- [23] J. Cai and Z.-S. Ye, “Contamination source identification: A bayesian framework integrating physical and statistical models,” *IEEE Transactions on Industrial Informatics*, vol. 17, no. 12, 2021.
- [24] J. Koch *et al.*, “Identification of contaminant source architectures—a statistical inversion that emulates multiphase physics in a computationally practicable manner,” *Water Resour. Res.*, vol. 52, no. 2, 2016.
- [25] J. Liang *et al.*, “Physics-informed data-driven models to predict surface runoff water quantity and quality in agricultural fields,” *Water*, vol. 11, no. 2, 2019.
- [26] R. Delabays *et al.*, “Locating the source of forced oscillations in transmission power grids,” *PRX Energy*, vol. 2, no. 2, 2023.
- [27] EPA, “EPA Stormwater Management Model (SWMM),” 2023. [Online]. Available: <https://www.epa.gov/water-research/storm-water-management-model-swmm>
- [28] A. Chio *et al.*, “GitHub Repository,” 2023. [Online]. Available: <https://github.com/andrewgchio/SWMMBackwardsInference>
- [29] M. Lubin *et al.*, “JuMP 1.0: Recent improvements to a modeling language for mathematical optimization,” *MPC*, 2023.
- [30] A. Wächter and L. T. Biegler, “On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming,” *Mathematical programming*, vol. 106, 2006.
- [31] UKRI, “Coin-hsl,” 2023. [Online]. Available: www.hsl.rl.ac.uk/ipopt