# Packaging your research problem: the case of *CityWaterBalance* for R

Laura Erban

EPA R User Group Workshop

August 12, 2019

# Introduction

Speaker:        hydrogeologist, earth systems scientist, postdoc

Location:        ORD/NHEERL/AED (Narragansett, RI)

AspiRation:      saving my future self

Disclaimer:      all opinions are my own

# Heuristics

**David Robinson**
@drob

**Follow** ∨

If you write the same code 3 times you should write a function

You should make an R package even for code that you don't plan to distribute. You'll find it is easier to keep track of your own personal R functions if they are in a package. And it's good to write documentation, even if it's just for your future self.

Source: https://kbroman.org/pkg_primer/

# Research problems

- o complexity (of subject, analysis, collaborators)

- o reproducibility

- o quantity / task-switching

# How much water flows through a city?



Icons made by Freepik from www.flaticon.com

# How much water flows through a city?
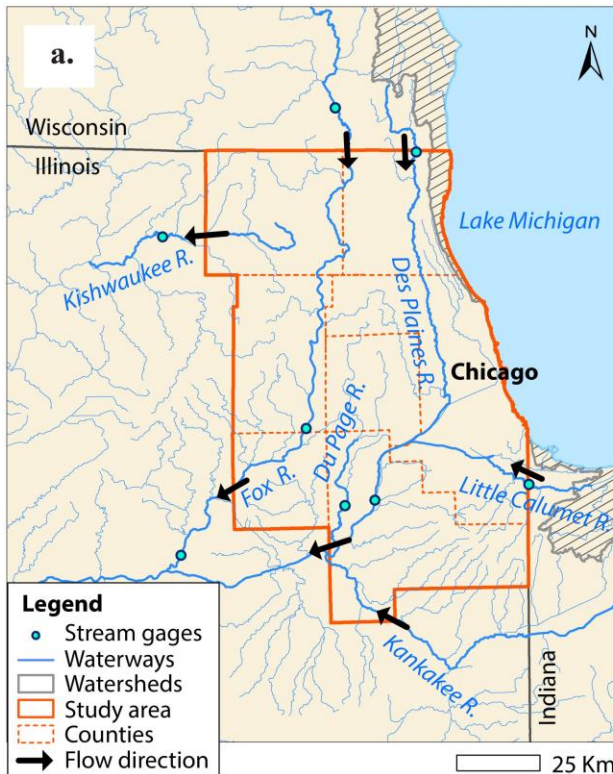


Icons made by Freepik from www.flaticon.com

- diverse sources / uses

- internal recycling / reuse

- scattered data

- unmeasured flows

# Water flows in greater Chicago*

*seven counties under regional planning agency:  Chicago Metropolitan Agency for Planning (CMAP)
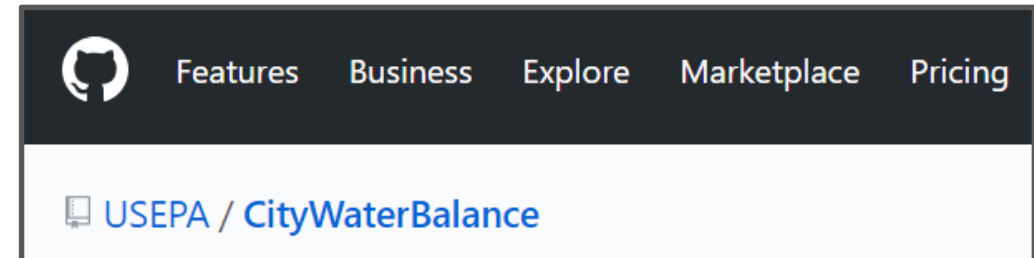
Study area

Developed land

Lake Michigan users



Erban et al., 2018 *Open Water Journal*

CMAP, 2017  *On to 2050: Water Resources*

# Packaging the workflow

*CityWaterBalance* for R

(available on CRAN and USEPA GitHub)

Design principles:

- open source, open access

- generalizable to any urban area

- automate data retrieval

- estimate unmeasured flows



Features    Business    Explore    Marketplace    Pricing

USEPA / **CityWaterBalance**

**Track flows of water through an urban system**

## CityWaterBalance

`build` `passing`

`CityWaterBalance` provides a reproducible workflow for studying an urban water system. The network of urban water flows and storages can be modeled and visualized. Any city may be modeled with preassembled data, but data for US cities can be gathered via web services using this package and dependencies.

## To install

Install the development version of CityWaterBalance from GitHub:

```
install.packages("devtools")
library(devtools)
install_github("USEPA/CityWaterBalance")
library(CityWaterBalance)
```

# Inside *CityWaterBalance*

Networked model:

- accounts for inflows, outflows, changes in storage

- based on visual math (right), coded in the core package function
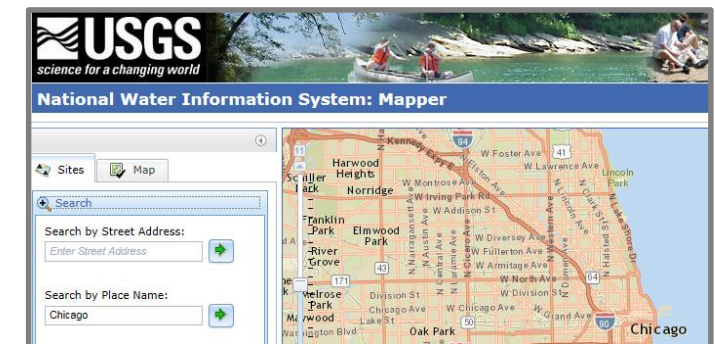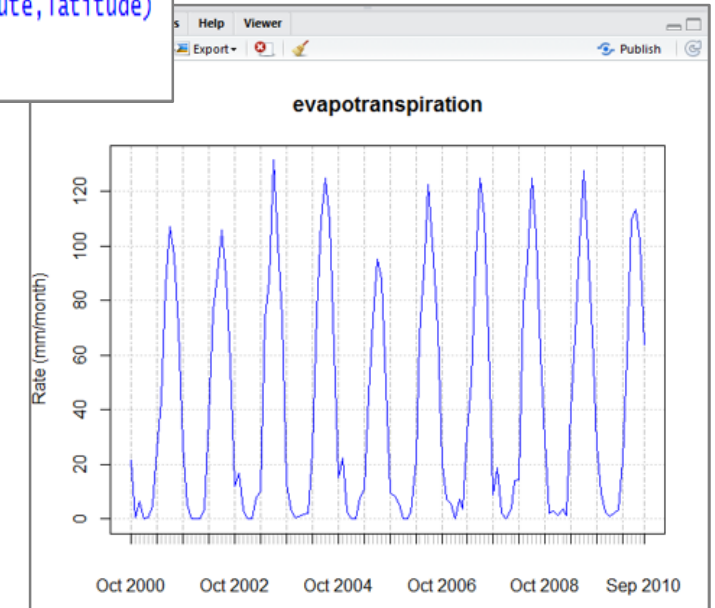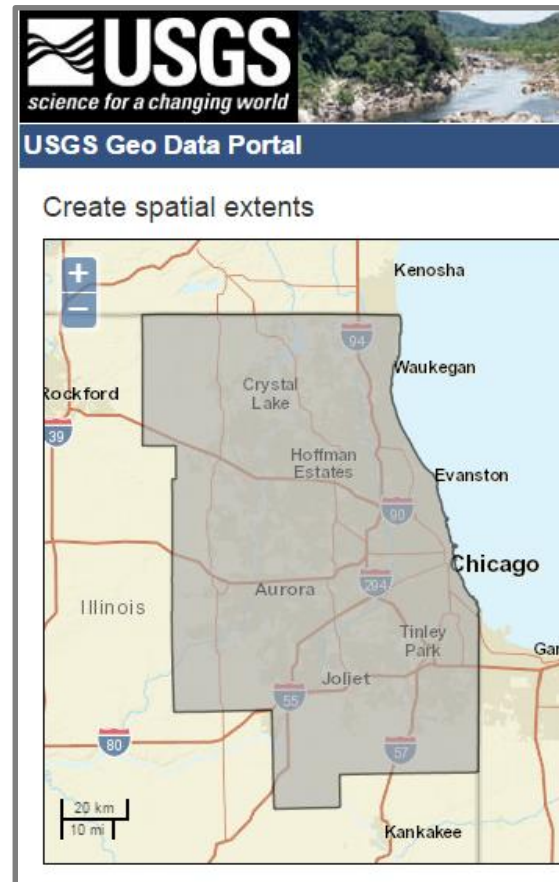


Urban water system

Key

◯ Global flows: precipitation (PRCP), streamflow in (INFLOW), imports (IMPORTS), evapotranspiration (ET), streamflow out (OUTFLOW)

⬡ Storages: surface water (SW), combined sewer system (CSS), shallow groundwater (SGW), deep groundwater (DGW)

▭ Producers: purification plants (PUR), industrial facilities (IND), wastewater treatment plants (WTP)

⬡ Consumers: potable use (POT), non-potable use (NPOT)

Erban et al., 2018 *Open Water Journal*

# Inside *CityWaterBalance*

Usage:

- functions retrieve data from web services

- help from USGS R packages *geoknife* and *dataRetrieval*

# Inside *CityWaterBalance*

Usage:

- add data that is not federated, or not served online

- package as placeholder for future accessible data

# Inside *CityWaterBalance*

Usage:

- estimate unmeasured flows by mass balance



Urban water system

Key
- Global flows: precipitation (PRCP), streamflow in (INFLOW), imports (IMPORTS), evapotranspiration (ET), streamflow out (OUTFLOW)
- Storages: surface water (SW), combined sewer system (CSS), shallow groundwater (SGW), deep groundwater (DGW)
- Producers: purification plants (PUR), industrial facilities (IND), wastewater treatment plants (WTP)
- Consumers: potable use (POT), non-potable use (NPOT)

dS/dt?

# *CityWaterBalance* output, from two perspectives

1) time

2) uncertainty
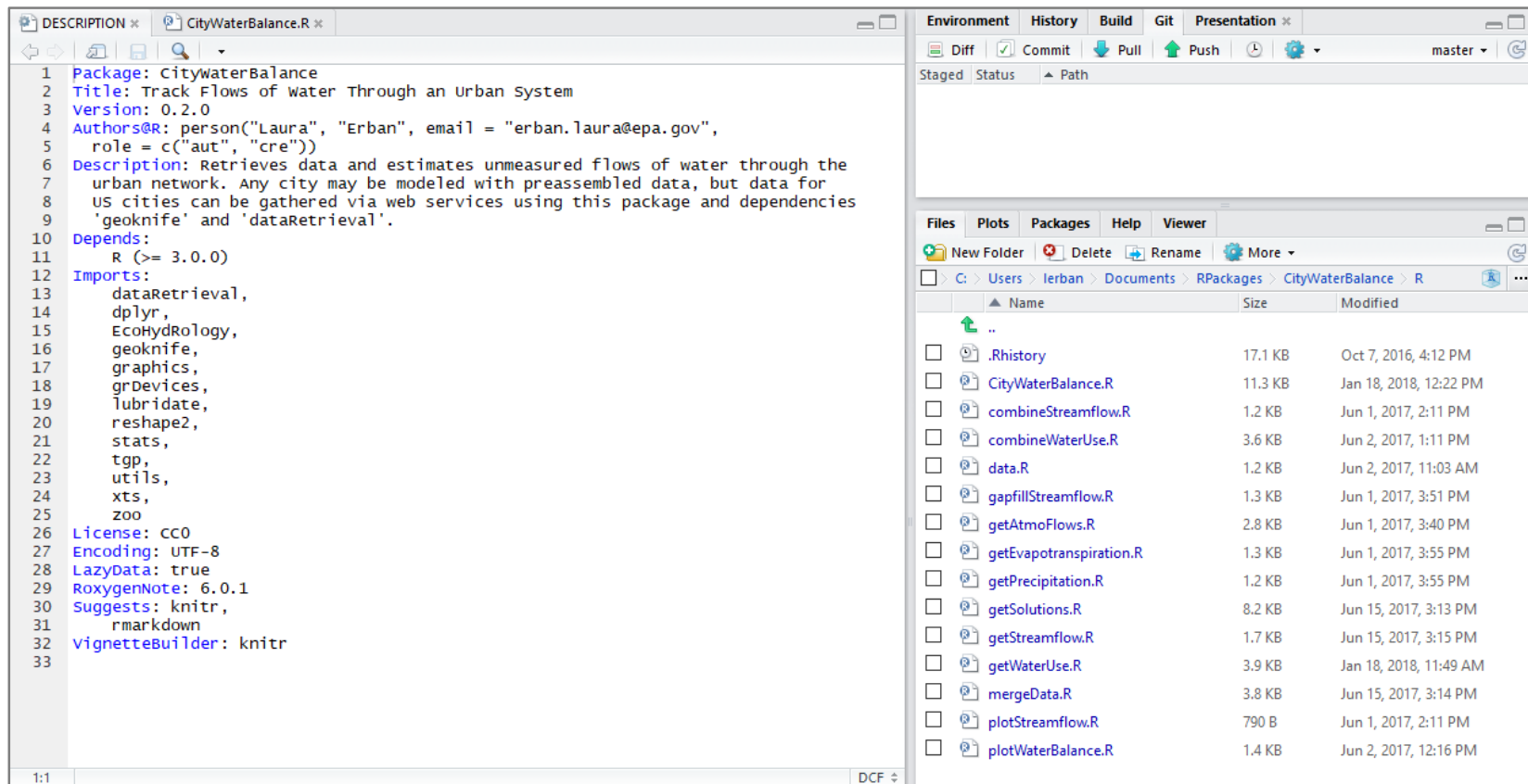
# A quantitative portrait



Average annual flows:  water years 2001-2010

# Packaging your research problem: generalizable lessons

# Writing a package can help you..

1. organize your functions

# Writing a package can help you..

2. document your functions in a common format

# Writing a package can help you..

## 3. share your work

**Install**
Usage overview
Usage examples

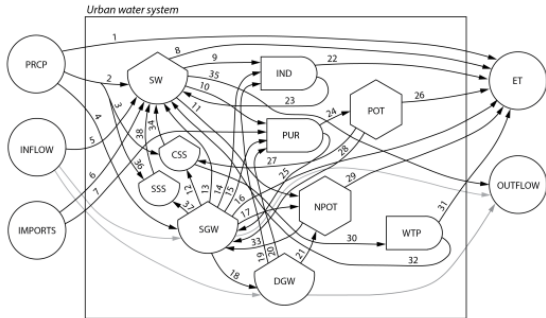# Introduction to CityWaterBalance

**Laura Erban**

**2018-02-13**

`CityWaterBalance` provides a reproducible workflow for studying urban water systems. Any system may be modeled with preassembled data, but data for US cities can be gathered via web services using this package and dependencies geoknife and dataRetrieval.

Install
**Usage overview**
Usage examples

# Usage overview

`CityWaterBalance` is based on a model of the urban water system, shown in the diagram below. This diagram specifies the network of water flows along with a mathematical solution for the changes in water storages (i.e., inflows - outflows) within the system.

*Urban water system*

**Key**
○ Global flows: precipitation (PRCP), streamflow in (INFLOW), imports (IMPORTS), evapotranspiration (ET), streamflow out (OUTFLOW)
○ Storages: surface water (SW), combined sewer system (CSS), separated sewer system (SSS), shallow and deep groundwater (SGW, DGW)
▢ Producers: purification plants (PUR), industrial facilities (IND), wastewater treatment plants (WTP)
○ Consumers: potable use (POT), non-potable use (NPOT)

Install
Usage overview
Usage examples
  Option 1: Input preassembled data
  **Option 2: Input data gathered from web services**
  Solve

# Option 2: Input data gathered from web services

`CityWaterBalance` has other functions that assemble data for the model. At this time, these functions access US-based web services.

## Specify spatial and temporal boundaries

Define an area of interest (AOI) and upload that geometry to the USGS Geo Data Portal (GDP). The GDP will give the geometry a name, which may start with 'upload:'. Here we use a geometry that is already available to the GDP in order to automate the example.

```
geometry <- 'sample:Counties'
attribute <- 'STATE'
value <- 'RI'
area <- 2707
start <- "2010-01-01"
end <- "2010-12-31"
```

## Get atmospheric data

```
latitude <- 41.5801
atm <- getAtmoFlows(start, end, geometry, attribute, value, latitude)
```

## Get streamflow data

Choose streamgages to evaluate total inflow and outflow for the AOI. NWIS mapper may be useful here.

```
ingages <- c("01112500")
outgages <-c("01113895","01114000","01117000","0111850
0")
```

## Considerations

- Lots of help available for writing a package for R

- Not every workflow needs to be a package

- R notebooks as one possible alternative

- Do make your steps traceable, in whatever format you choose

# Thank you!

erban.laura@epa.gov

@leerban

 lerban