# Understanding the relationship between indicators & tracers and vapor intrusion

## Dynamic multivariate time series regressions

Riley Mulhern, RTI International

Chris Lutes, Jacobs

Chase Holton, GSI Environmental
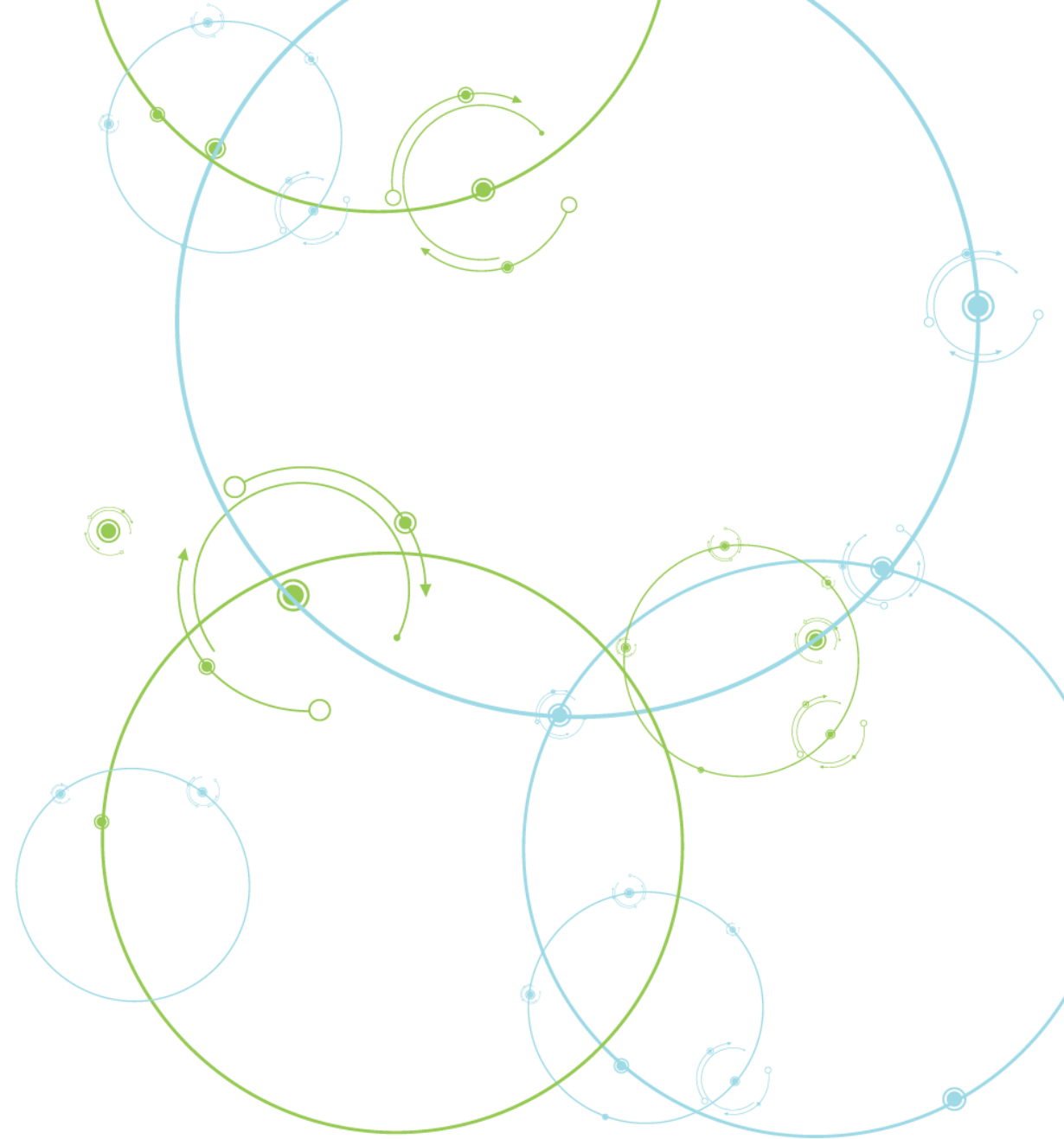
AJ Kondash, RTI International

**AEHS Annual Conference | U.S. EPA "State of VI Science" Workshop | March 20-23, 2023**

# Presentation Overview

o Modeling Aims & Approach

o Data Wrangling & Preprocessing

o Regression Development

o Results

o Conclusions & Next Steps

# Modeling Aims & Approach

# Challenges with using indicators and tracers to assess VI

- o Radon
- o Climatic conditions
- o Building conditions

?

→

Indoor air VOC concentrations

**Long term goals:**

- o Guide sampling decisions
- o Early warnings
- o Mitigate exposures
- o Soil Gas Safe Communities

# Primary aim is to gain insights into the relationship between indicators/ tracer and VI (*not prediction of future VI*)

o Radon
o Climatic conditions
o Building conditions

?

→

Need better understanding of the relationship between indicators/tracers and VI **across sites** and **over time** to make **generalizable recommendations**.

Indoor air VOC concentrations

# Challenges with using indicators and tracers to assess VI

o Radon
o Climatic conditions
o Building conditions

?

→

Indoor air VOC concentrations

**Past work:**

o Single variate time series analyses

o Multivariate time series analysis for a single site

# Specifying regression predictor variables and outcome variable

**Predictor variables**

- Radon (R, pCi/L)
- Differential temperature (ΔT, °C)
  - indoor-outdoor
- Differential pressure (ΔP, Pa)
  - positive indicates higher indoor pressure
  - Most sites: indoor-subslab
  - Indianapolis 422 first floor: basement-first floor

**Outcome variable**

Indoor air VOC concentrations (C, $\mu g/m^3$)

# Temporal data sets from three different VI sites across the country

## Virginia Site A



o Coastal military site, 120,000 ft²
o Brick with poured concrete slab 6-8 in. thick
o Separate HVAC zones
o **3 sampling sites** (Office, Supply room, Women's Restroom)
o ~19 months of data
o **Trichloroethylene (TCE)** from historical releases of chlorinated solvents near the site

## Sun Devil Manor



o Layton, UT
o Modern suburban residential
o Split level, ground floor sampling location
o **1 sampling site**
o ~17 months of data
o Trichloroethylene (TCE)

## Indianapolis 422



o Indianapolis, IN
o Constructed ~1915
o Duplex
o Wood frame, brick foundation, concrete basement floor
o **2 sampling sites** (Basement, First Floor)
o ~4 months of data
o **Perchloroethylene (PCE)** from historical dry cleaning and adjacent businesses

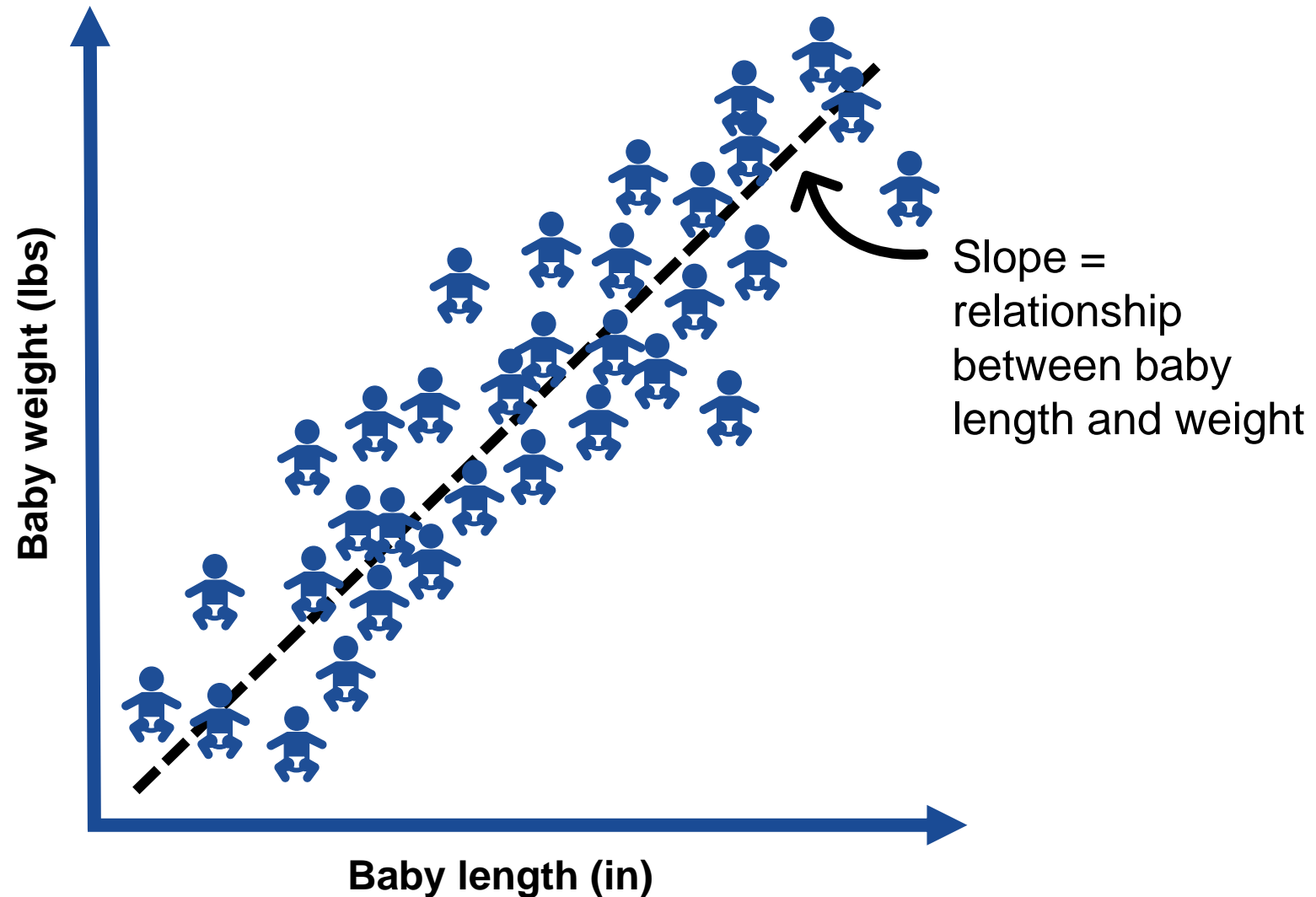# Ideal data for traditional linear regression: a simple example

o Assumes all observations are **independent** (not related to each other)

- Ex. What is the relationship between **length** and **weight** among babies less than 1 year old in the U.S.?

- Hypothesis: As babies grow longer, their weight increases.

- Ideal theoretical data set to test this hypothesis using linear regression: single measurements of many individual babies' weights and lengths → **a little bit of data from a lot of different babies**

Length

Weight

# **Ideal data** for traditional linear regression: a simple example

# Ideal data for traditional linear regression: a simple example



Slope = relationship between baby length and weight

Baby weight (lbs)

Baby length (in)

# Ideal data for traditional linear regression: a simple example

**Fitted vs. Residuals plot**

✓ **Accurate estimates** of regression coefficients (true relationship between the variables)
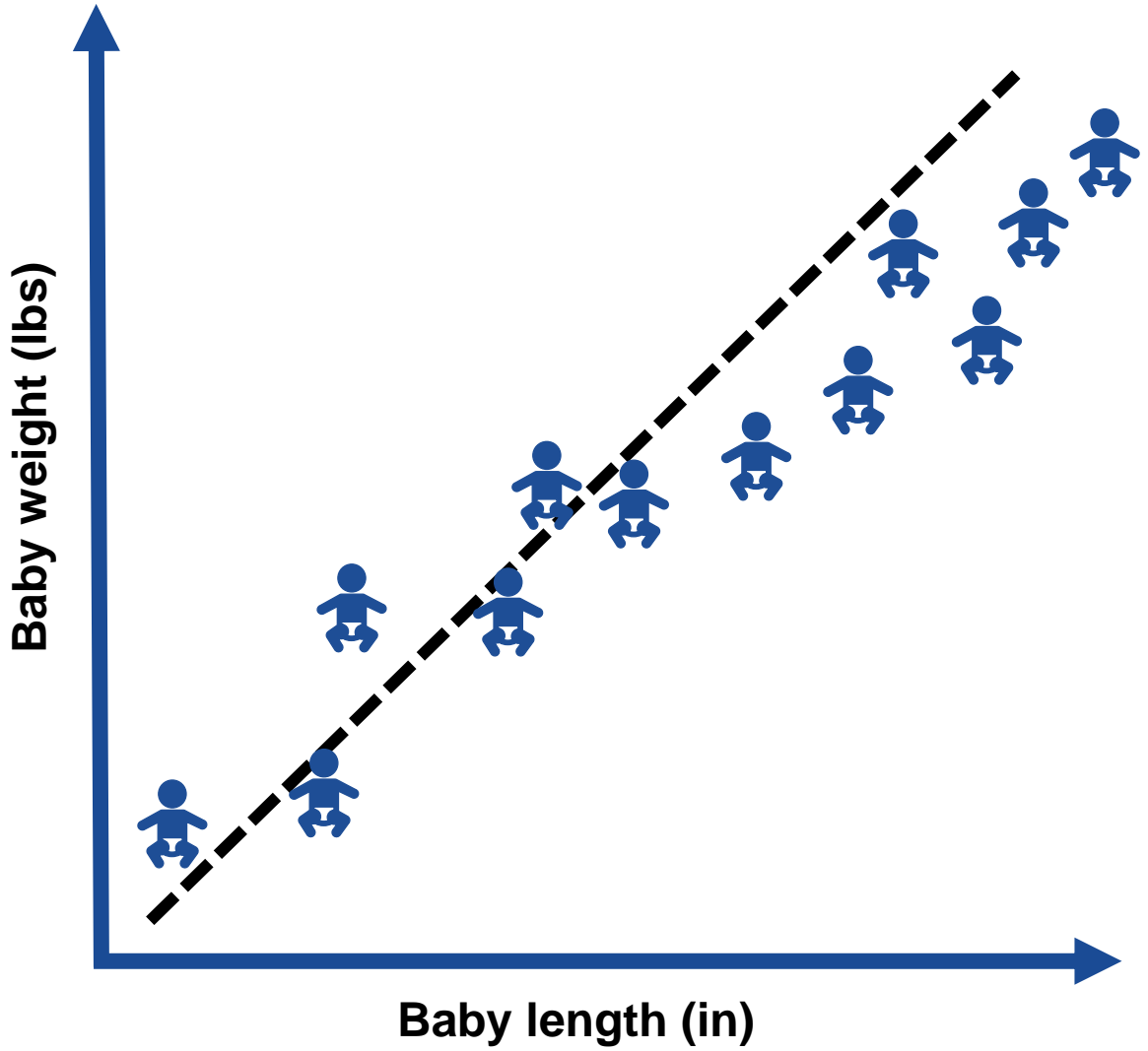
✓ **High confidence** in the significance of the relationship

Regression residual

Predicted weight (lbs)

This line indicates a perfect fit

# Temporal data set: challenges to using traditional regression

o Actual data set: multiple measurements from a single baby at different points in time → **a lot of data from only one baby**

o Hypothesis: As babies get older, their weight increase.

o Time series data are **not independent** → each measurement is related to measurements in the past

Riley's baby

# Temporal data set: challenges to using traditional regression

# Temporal data set: challenges to using traditional regression

✕ **Autocorrelation**
Leads to inaccurate estimates of the regression coefficients because we ignore important information in the data set.

# Temporal data set: challenges to using traditional regression

✗ **Spurious regression**
Could conclude that a significant relationship exists between variables when it doesn't actually exist
- Air transport in Australia related to rice production in New Guinea?

# Temporal data set: challenges to using traditional regression

**Fitted vs. Residuals plot**



✗ **Systematic bias**
Measurement errors over time can lead to systematic bias in predictions

Systematic pattern in residuals for higher predicted weights

**Predicted weight (lbs)**

# Actual data set: time series data from three sites

o Many observations of VOC concentrations and predictor variables at sequential points in time from three sites across the country

# Actual data set: time series data from three sites

o Many observations of VOC concentrations and predictor variables at sequential points in time from three sites across the country

# **Actual data set:** time series data from three sites

**Risk of:**

o Autocorrelation and inaccurate estimates

o Spurious regression results

o Systematic bias

Site 1

Site 2

Site 3

# **Actual data set:** time series data from three sites

**Risk of:**

o Autocorrelation and inaccurate estimates

o Spurious regression results

o Systematic bias



*Cannot know the <u>true</u> relationship between the predictor and the outcome variable using traditional linear regression for time series data.*

# Dynamic Time Series Regression

Refers to the "dynamic" (non-static) nature of time series data where past observations influence current and future values

**Regression approach**

- ○ **Outcome variable:** Indoor air VOC concentration (C, μg/m³)
- ○ **Predictor variables:**
  - Indoor air radon concentration (R, pCi/L)
  - Differential temperature (ΔT, °C)
  - Differential pressure (ΔP, Pa)

**Traditional linear regression:**

$$\mathrm{C}_t = \boldsymbol{\beta_0} + \boldsymbol{\beta_1}\mathrm{R}_t + \boldsymbol{\beta_2}\Delta T_t + \boldsymbol{\beta_3}\Delta P_t + \varepsilon_t$$

**Dynamic time series regression:**

$$\mathrm{C}_t = \boldsymbol{\beta_0} + \boldsymbol{\beta_1}\mathrm{R}_t + \boldsymbol{\beta_2}\Delta T_t + \boldsymbol{\beta_3}\Delta P_t + \eta_t$$

**"ARIMA errors"**
Allows the residuals of the regression to follow an autoregressive integrated moving average (ARIMA) regression

# Dynamic Time Series Regression

**Traditional linear regression:**

$$C_t = \boldsymbol{\beta_0} + \boldsymbol{\beta_1} R_t + \boldsymbol{\beta_2} \Delta T_t + \boldsymbol{\beta_3} \Delta P_t + \varepsilon_t$$

**Dynamic time series regression:**

$$C_t = \boldsymbol{\beta_0} + \boldsymbol{\beta_1} R_t + \boldsymbol{\beta_2} \Delta T_t + \boldsymbol{\beta_3} \Delta P_t + \eta_t$$

**Today**
Try to use todays' differential temperature, differential pressure, and radon to predict today's VOC concentration.

C

Time

**Yesterday**

**Today**
Try to use todays' differential temperature, differential pressure, and radon **and how wrong my prediction was from yesterday** to predict today's VOC concentration.

*Allows the regression to "control" for the temporal nature of the data to accurately estimate the true relationship between the predictor variables and the outcome variable.*

C

Time

# Data Wrangling and Preprocessing

### Virginia Site A



### Sun Devil Manor



### Indianapolis 422

# Data Summary

**SITES**

**SAMPLE LOCATIONS**

**TIME PERIODS**

**AVERAGING TIMES**

**Virginia Site A**

Supply Room → 5/2017 – 1/2021

Women's Restroom → 11/2019 – 1/2021

Office → 7/2019 – 1/2020

**Sun Devil Manor**

Ground level → 1/2011 – 6/2012

**Indianapolis 422**

Basement → 8/2011 – 10/2011

Basement → 12/2011 – 2/2012

First Floor → 8/2011 – 10/2011

First Floor → 12/2011 – 2/2012

- 6 hour
- 24 hour
- Weekly

# Virginia Site A: **Supply Room**



Identify region of dataset with overlapping regressor and response variables

**tsibble**

Create time series object

**fable**

Interpolate missing data intervals (auto ARIMA)

# Virginia Site A: **Women's Restroom**



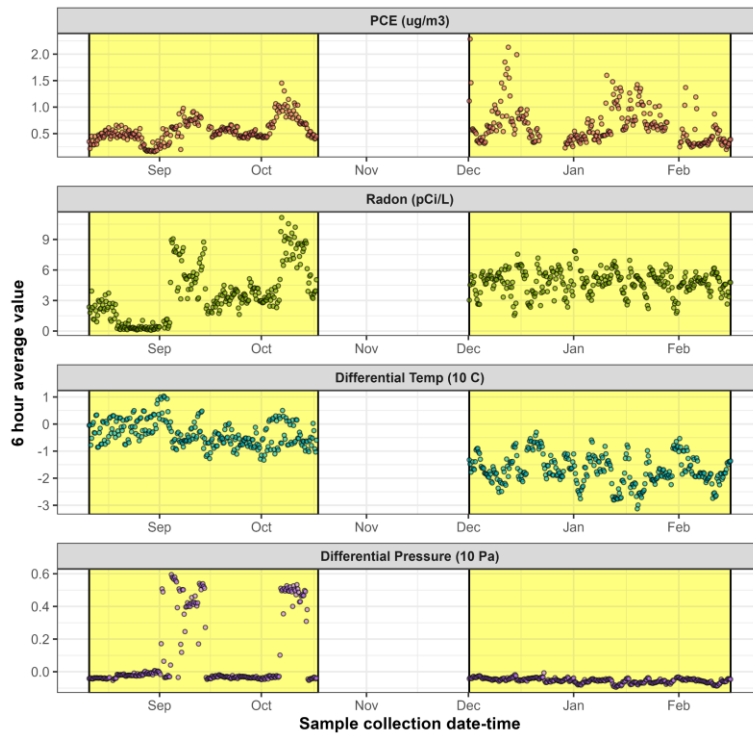Identify region of dataset with overlapping regressor and response variables → Create time series object → Interpolate missing data intervals (auto ARIMA)
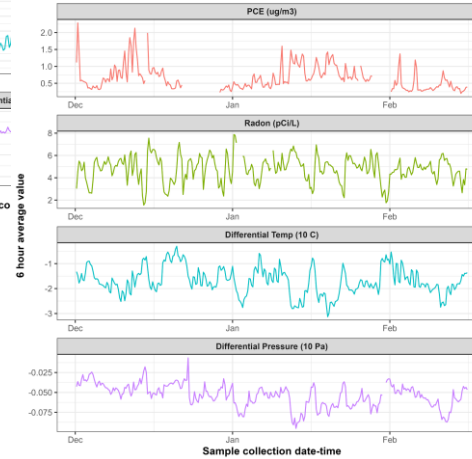
# Virginia Site A: **Office**



Identify region of dataset with overlapping regressor and response variables → Create time series object → Interpolate missing data intervals (auto ARIMA)

# Sun Devil Manor



Identify region of dataset with overlapping regressor and response variables → Create time series object → Interpolate missing data intervals (auto ARIMA)

# Indianapolis 422 **Basement**



Identify region of dataset with overlapping regressor and response variables

→

Create time series object

→

Interpolate missing data intervals (auto ARIMA)

# Indianapolis 422 **First Floor**



Identify region of dataset with overlapping regressor and response variables → Create time series object → Interpolate missing data intervals (auto ARIMA)

# Regression Development

# 8 time series, 24 multivariate regressions



- 3 **sites**
  - 6 total **sample locations**
    - 8 total **time periods**

- Each time series modeled using **three different averaging times**
  - 6 hour
  - 24 hour
  - Weekly

- All regression combinations tested:
  - 24 **complete models** with all predictors
  - 72 **"leave-one-out"** models
  - 72 **single variate models**
  - 168 total models

# Develop a separate regression for each time series and averaging time

## Sun Devil Manor

**6 hour**



Outcome variable

Predictor variables

**24 hour**



Outcome variable

Predictor variables

**Weekly**



Outcome variable

Predictor variables

$$C_t = \boldsymbol{\beta_0} + \boldsymbol{\beta_1} R_t + \boldsymbol{\beta_2} \Delta T_t + \boldsymbol{\beta_3} \Delta P_t + \eta_t$$

fable

feasts

# Results

**Virginia Site A**



**Sun Devil Manor**



**Indianapolis 422**

# Virginia Site A: **Supply Room**
## Complete Model Fit



**6 hour**

Fitted vs actuals

$R = 0.89, p < 2.2e\text{-}16$

**24 hour**

Fitted vs actuals

$R = 0.89, p < 2.2e\text{-}16$

**Weekly**

Fitted vs actuals

$R = 0.95, p < 2.2e\text{-}16$

# Virginia Site A: **Supply Room**
## Complete Regression Diagnostics



**Weekly**

Residuals should resemble a white noise series

Minimal autocorrelation

Residuals should represent a normal distribution

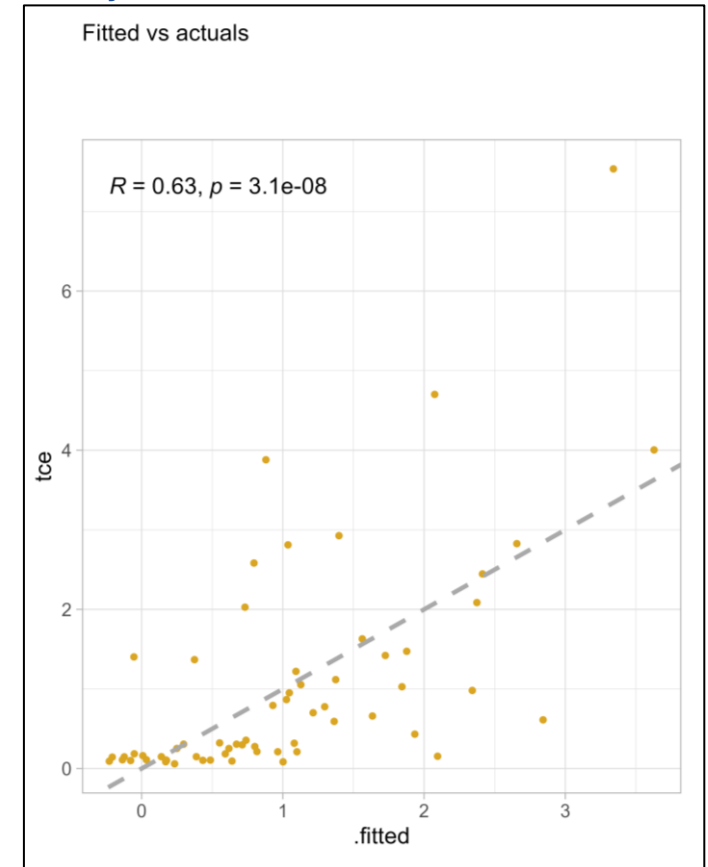# Virginia Site A: **Women's Restroom**
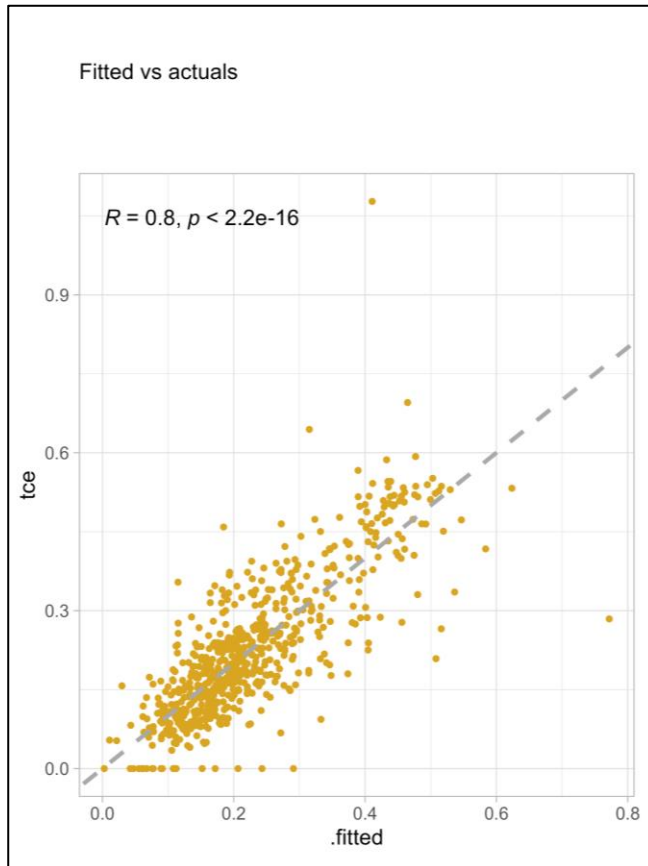## Complete Model Fit

**6 hour**

**24 hour**

**Weekly**

# Virginia Site A: **Office**
## Complete Model Fit



**6 hour**

Fitted vs actuals

$R = 0.8$, $p < 2.2$e-16

**24 hour**

Fitted vs actuals

$R = 0.82$, $p < 2.2$e-16

**Weekly**

Fitted vs actuals

$R = 0.85$, $p = 2.7$e-08

# Sun Devil Manor
## Complete Model Fit

**6 hour**

Fitted vs actuals

$R = 0.76, p < 2.2e\text{-}16$

**24 hour**

Fitted vs actuals

$R = 0.8, p < 2.2e\text{-}16$

**Weekly**

Fitted vs actuals

$R = 0.79, p < 2.2e\text{-}16$

## Complete Model Fit

**6 hour**



**24 hour**



**Weekly**

# Indianapolis 422 **Basement**: Dec-Feb
## Complete Model Fit

**6 hour**



Fitted vs actuals

$R = 0.89$, $p < 2.2e\text{-}16$

**24 hour**



Fitted vs actuals

$R = 0.88$, $p < 2.2e\text{-}16$

**Weekly**



Fitted vs actuals

$R = 0.53$, $p = 0.076$

# Indianapolis 422 **First Floor**: Aug-Oct
## Complete Model Fit

**6 hour**

Fitted vs actuals

$R = 0.94, p < 2.2\text{e-}16$

pce / .fitted

**24 hour**

Fitted vs actuals

$R = 0.94, p < 2.2\text{e-}16$

pce / .fitted

**Weekly**

Fitted vs actuals

$R = 0.81, p = 0.0022$

pce / .fitted

# Indianapolis 422 **First Floor**: Dec-Feb
## Complete Model Fit

**6 hour**



Fitted vs actuals

R = 0.83, p < 2.2e-16

**24 hour**



Fitted vs actuals

R = 0.74, p = 9.5e-15

**Weekly**



Fitted vs actuals

R = 0.54, p = 0.067

$$C_t = \boldsymbol{\beta_0} + \boldsymbol{\beta_1}R_t + \boldsymbol{\beta_2}\Delta T_t + \boldsymbol{\beta_3}\Delta P_t + \eta_t$$

$$C_t = \boldsymbol{\beta_0} + \boldsymbol{\beta_1}R_t + \boldsymbol{\beta_2}\Delta T_t + \boldsymbol{\beta_3}\Delta P_t + \eta_t$$

# Six hourly radon concentrations significant across all sites, sample locations, and time periods



- Across all sites, locations, and time periods, **every 1 pCi/L increase** in 6 hourly average radon concentration results in a **0.04-5.36 µg/m³ increase** in 6 hourly VOC concentrations (median=0.14 µg/m³).

- Could use this relationship to determine when a household may be at risk of exceeding VOC screening level based on indoor air radon.

# Radon **less reliable** as a tracer at **longer averaging times**

# Significance of other indicators and tracers **site specific** and **dependent on** sample averaging time



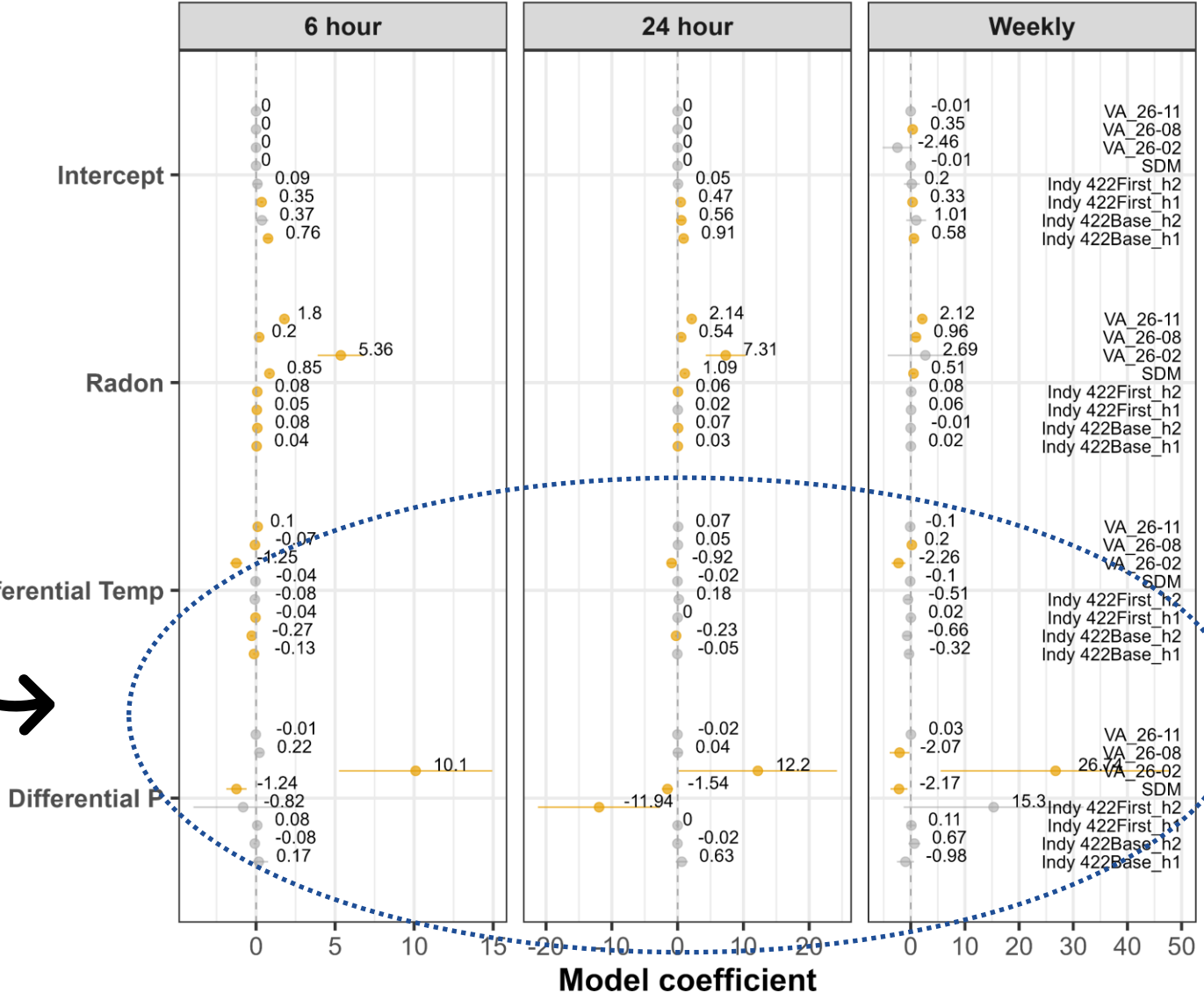Differential pressure and temperature dependent on the **sample location** within the building and the **sampling period**.

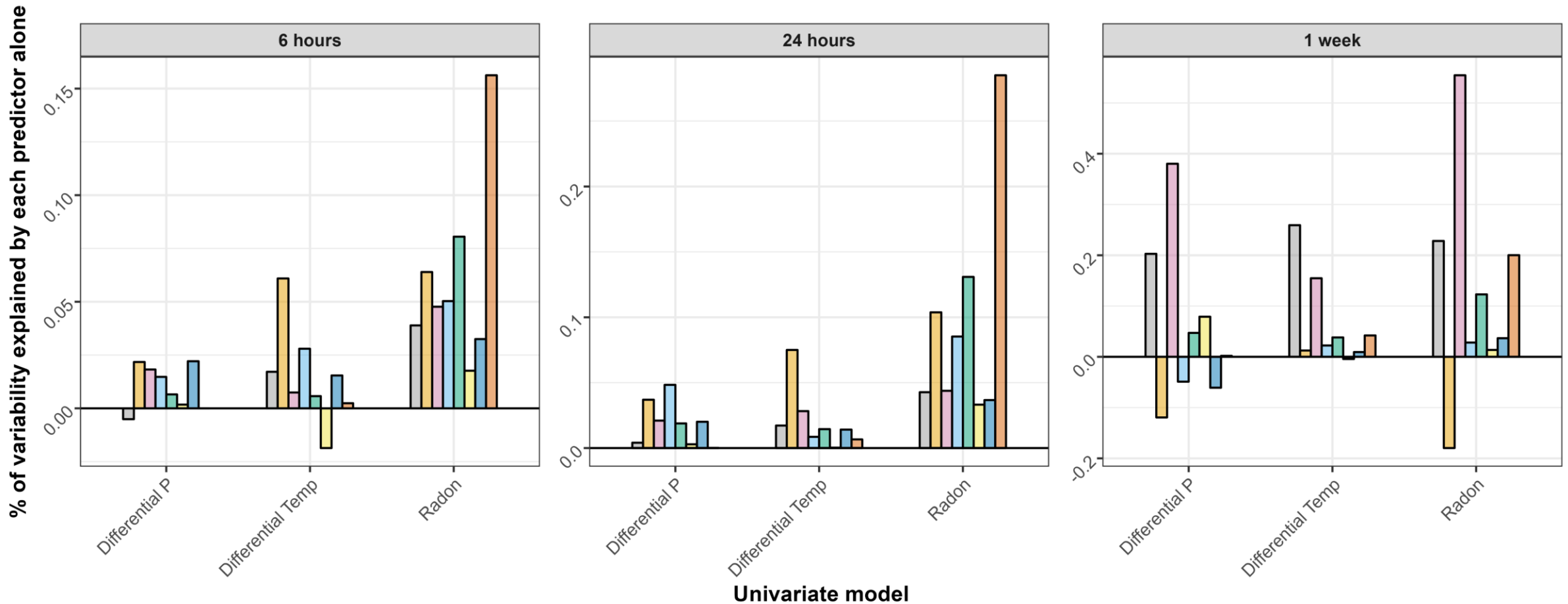# **Variable importance** – How much does the fit improve with a univariate model compared to a *null model*?

Greater increase = greater importance



**Site**
- Indy 422Base_h1
- Indy 422Base_h2
- Indy 422First_h1
- Indy 422First_h2
- SDM
- VA_26-02
- VA_26-08
- VA_26-11

# **Variable importance** – How much does the fit improve with a univariate model compared to a *null model*?

Greater increase = greater importance
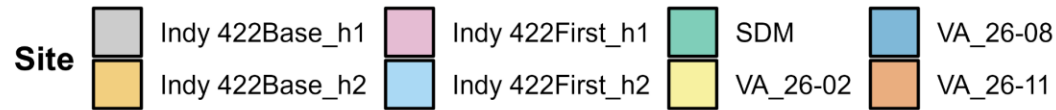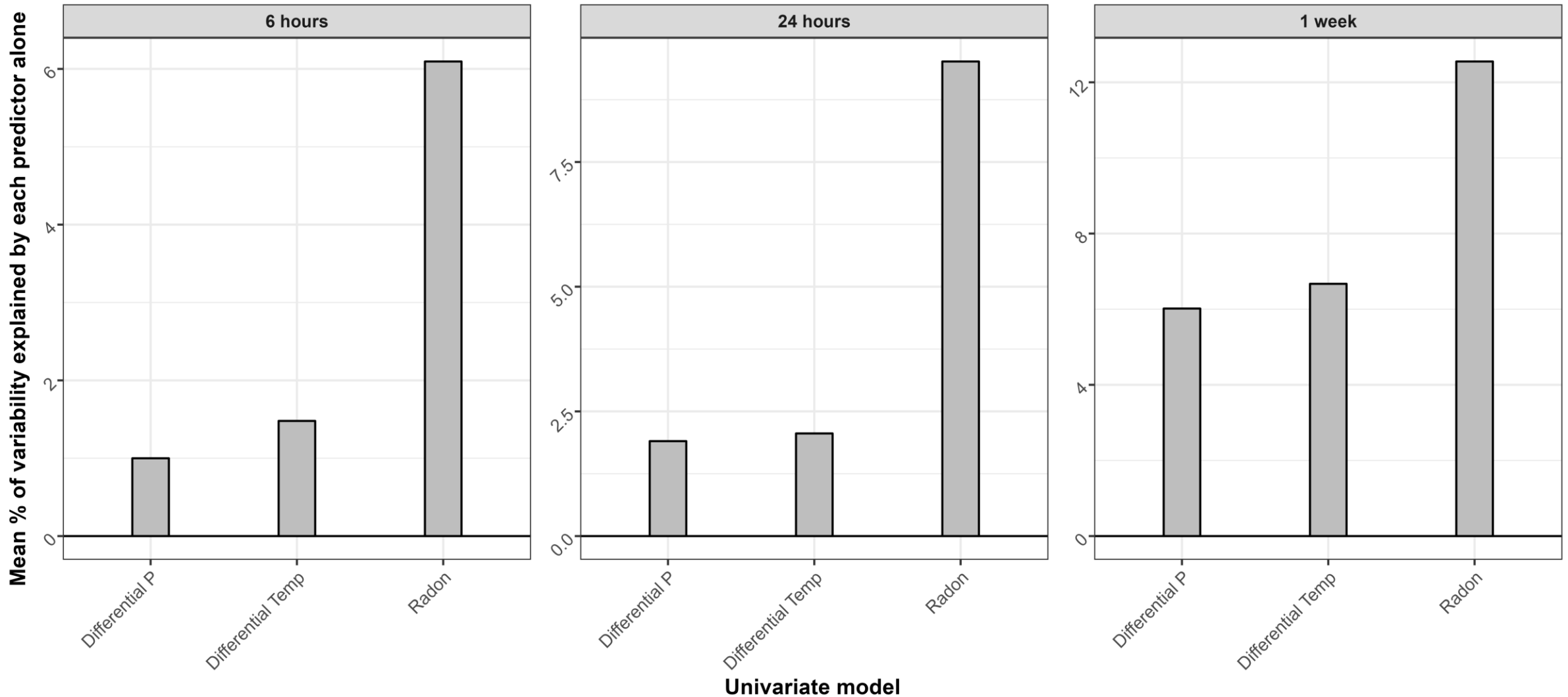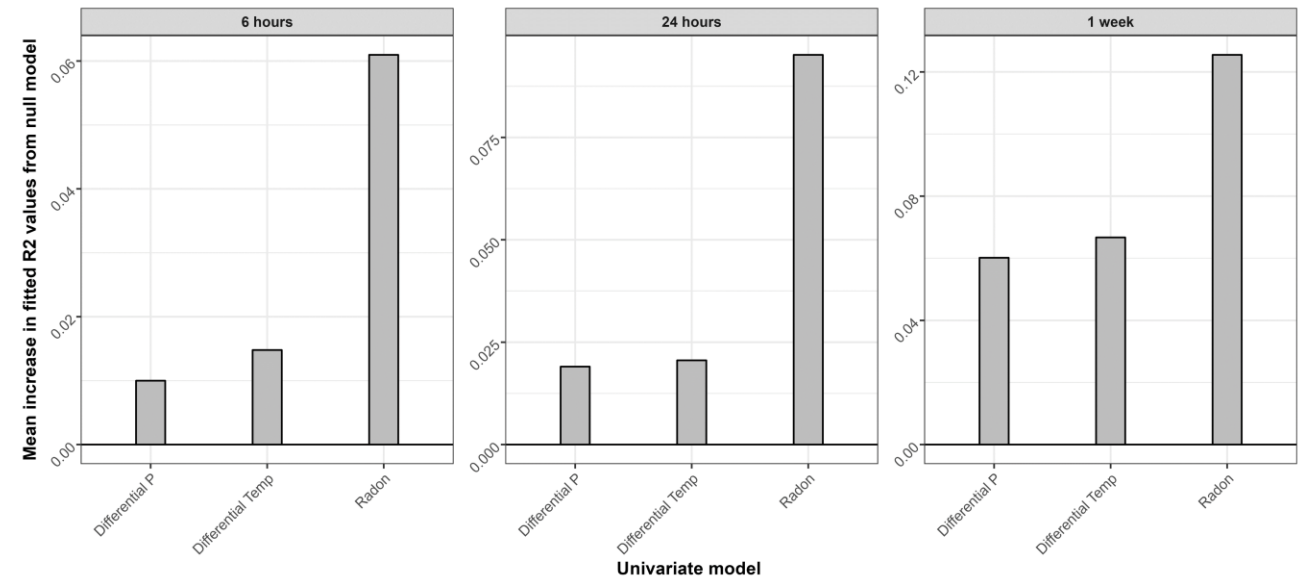
# Variable importance – How much does the fit improve with a univariate model compared to a *null model*?

Greater increase in $R^2$ = greater importance

- On average, **radon independently explains an additional 6-12% of the variability** in VOC concentrations beyond the temporal autocorrelation alone.

- **Differential temperature and differential pressure each explain only 1-7% of the variability** in VOC concentrations beyond the temporal autocorrelation alone.

# Conclusions & Next steps

o **6-hour radon concentration is the most reliable tracer** for indoor air VOC concentrations across sites after controlling for temporal autocorrelation

- Relationship between Δ6-hour radon and Δ6-hour VOC varies by site
- Baseline relationship could be characterized to develop general recommendations for different areas
- Simple in-home radon detectors may be the best, most immediate indicator of increased VOC exposure risk in VI areas

o **Other indicators are largely site specific** and vary according to location within the building and time of year

o **Future analysis:** Include additional covariates, include additional sites, forecasting future VI

# Thank you

Contact: Riley Mulhern | email: rmulhern@rti.org