# Data Science in Production: Making Data Science Accessible for Domain Experts and Creating Real Business Value

Kenneth A. Longo, Ph.D.
2020-11-19

KNIME Fall Summit 2020

# A problem and a solution

- **Problem**

  - How can we quickly and reliably develop, test and deploy software that uses modern tools of data science and answers key business questions?

  - Is there a platform that provides end-to-end software development tools for data science backend, data connectivity, UI, UX?
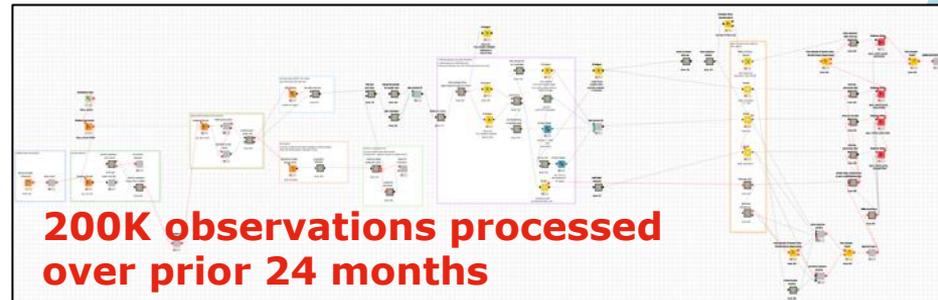
- **Solution**

  - KNIME Analytics Platform, Server & WebPortal

  - Provides access to 1) robust data science tools and 2) a means for productionizing solutions for broad use, delivery via browser.

WAVE™
LIFE SCIENCES

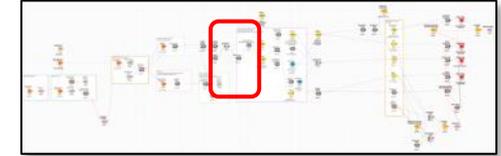# KNIME Server pipeline & Webportal GUI for qPCR

The **qPCR workflow** for **mRNA measurement** is a guided analytic process for scientists that analyses 384-well bioassay data, combining:

- **Ready-made features of KNIME**
  - **Components** with interactive JS graphs and tables for **guided metadata selection** & examining data
  - **Context properties** extraction (who, what, when...)
  - **Controlled metadata** = flow variables = protocol def.

- **Advanced statistical features (R)**
  - Linear mixed effect (LME) models
  - Robust fitting, outlier handling
  - R:ggplot2 graphics

- **AWS & Db connectivity**
  - Metadata Selection, compound validation & Db write



**200K observations processed over prior 24 months**

# Perform interactive visual examination of the data (1)



- Assess plate heatmaps for:
  - Obvious data deformations
  - Data loss
  - Outliers
  - Systematic column- or row-based effects
- Users can request assistance or insights by sending DS the job link

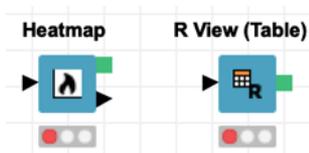# Perform interactive visual examination of the data (2)
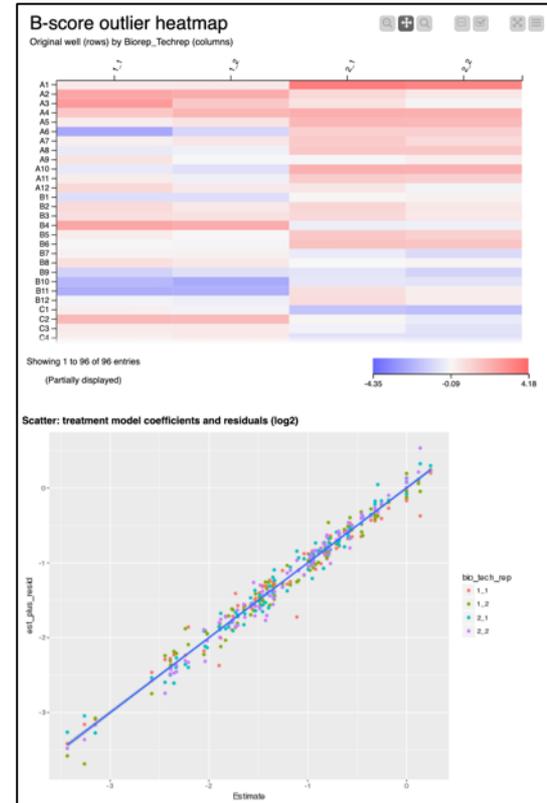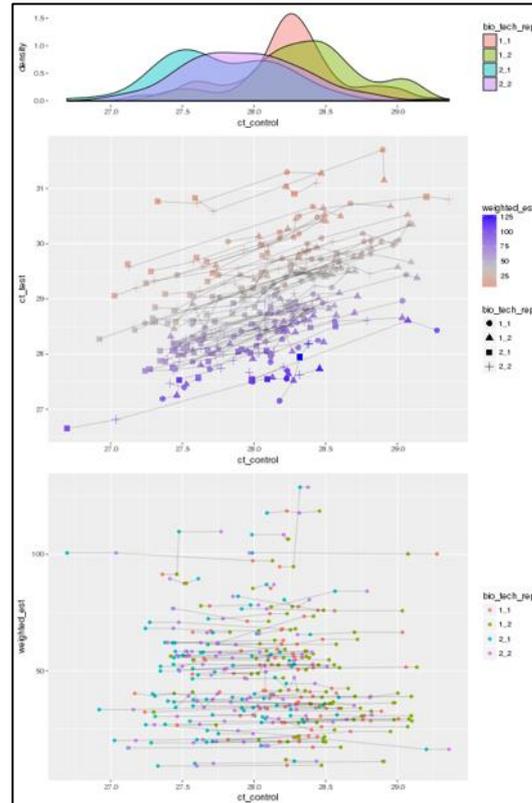


- Interactively view scatterplot of test and control gene raw data and its connected metadata in tabular form using KNIME js.nodes

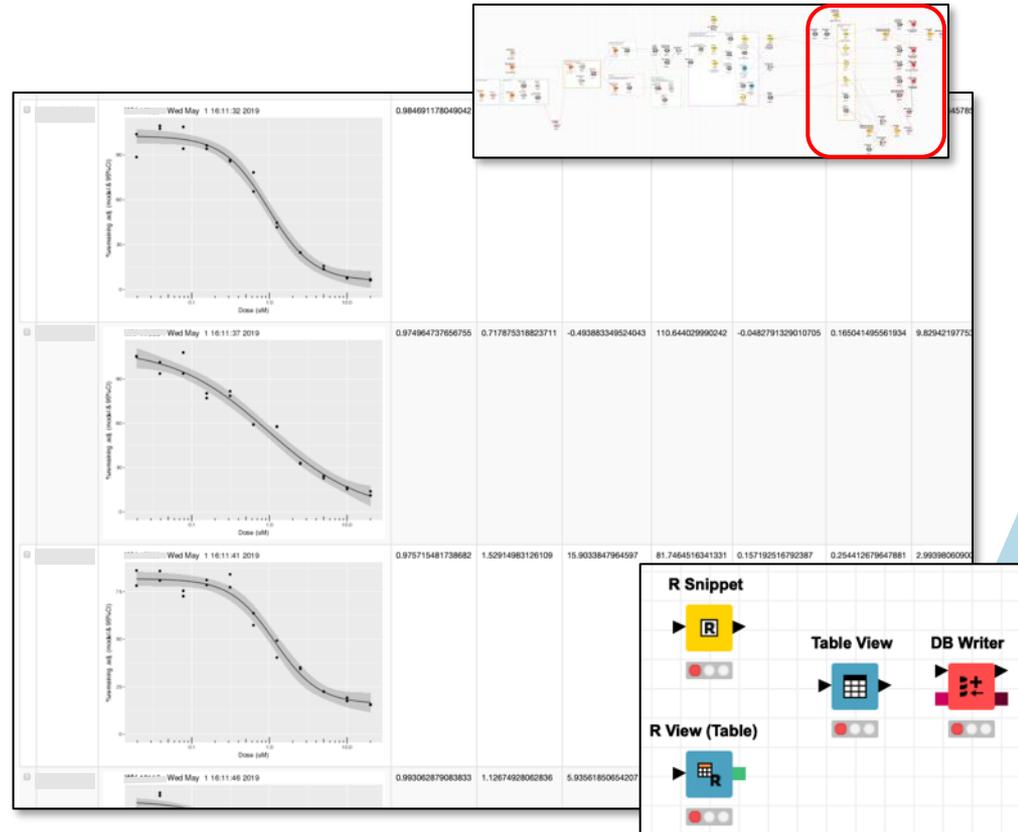- Generate initial impressions of data quality and likely outliers, reasons behind, etc

# Review model fit & outlier adjustments

- Examine how the rLME model dealt with outliers & impacted the final parameters

- Differentiate between technical and biological variation

- Uses combination of:
  - R:ggplot2 graphics
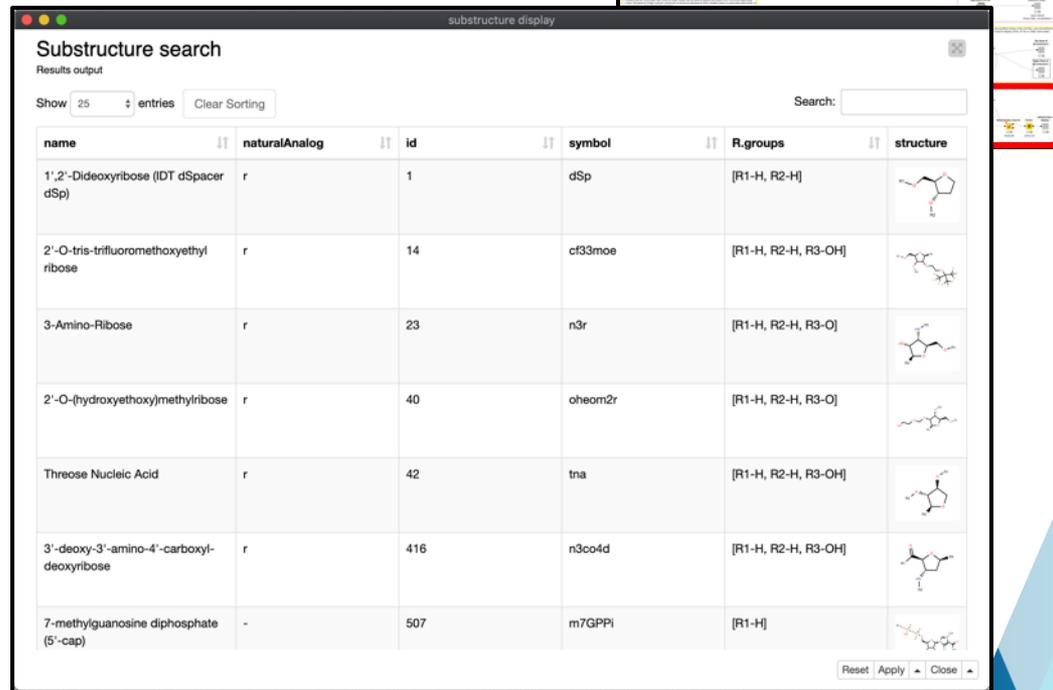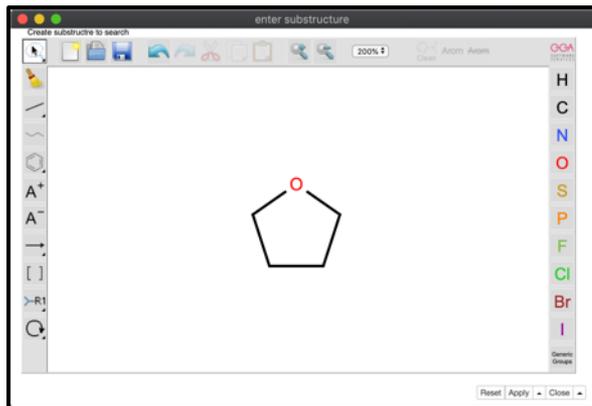  - KNIME js.nodes for interactive visualization of outliers
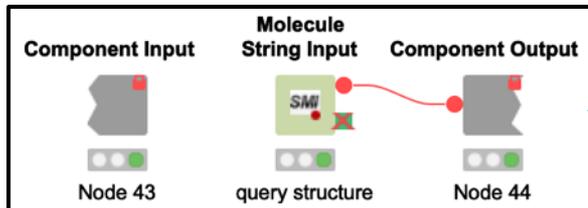
# Dose-response Kinetics and Database Commit

- 4-parameter dose-response kinetics and graphing using *R:drc* and *R:ggplot2*

- If-switch on > n-independent doses

- Review and QC using JS Table View
  - Graphs
  - Parameters
  - Confidence intervals
  - Goodness-of-fit stats

- Commit to AWS-resident Db
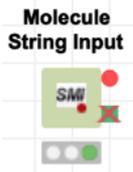
# Cheminformatics: substructure search

# Chemical display, categorization, selection

- **Composed** our own monomer library
  - **Used exisiting** minimal set (providing majority coverage)
  - **Supplemented** with proprietary monomers created using the *Molecule String Input* node, complete with R-groups

- Data Science/Chemistry **collaboration** framework!

# Bioinformatics: gene-species alignment network



Bowtie gene-species
alignment table w/
mismatch count

Species list

- Feeds from the bowtie alignment table; users select species of interest

- Webportal view using Network Viewer JS node

- Establishes relationship of sequence with multiple species, inversely
  edge-weighted by mismatch count, including off-target hypotheses