# WAVE®
## LIFE SCIENCES

Blending KNIME Data ETL and Analytics Capabilities with an Enterprise LIMS, ELN and Cheminformatics Platform

Ken Longo

June 1, 2022

# The Big Picture

Wave Life Sciences is a genetic medicines company focused on delivering life-changing treatments for people battling devastating diseases.

Early in the drug discovery process we use a combination of well-established enterprise software solutions for molecule registration, cheminformatics and data acquisition.

We use KNIME to interface with these platforms and perform a range of critical data ETL, modeling and analytics services.

WAVE
LIFE SCIENCES
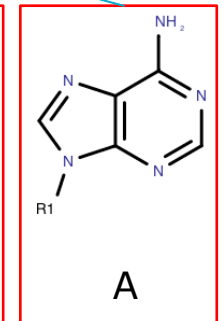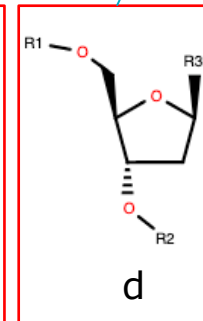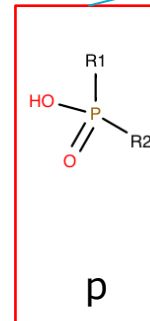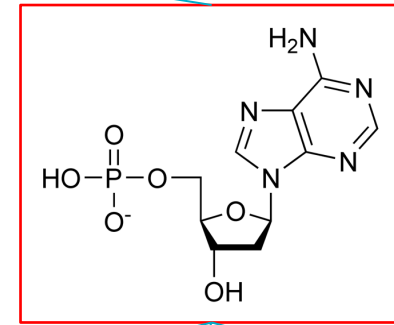
# A Story of Growth & Change at Wave…

- Grew from a small biotechnology start-up of <10 people into a 225+ person company

- Sponsors clinical trials worldwide in the areas of Huntington's Disease, Frontotemporal Dementia, Amyotrophic Lateral Sclerosis and Muscular Dystrophy

- Has a pipeline of several early-stage drug discovery programs that require an informatics infrastructure

- As we evolved our systems of organizing chemical entities and the data connected to them, our use of KNIME grew and changed also



WAVE®
LIFE SCIENCES

# Introduction to Wave Life Sciences

- Wave Life Sciences is a **genetic medicines** company focused on delivering life-changing treatments for people battling devastating diseases.

- Requires solutions for several domains:
  - **Bioinformatics**
  - **Cheminformatics**
  - **Data processing/ETL** – assay data
  - **Analysis** (graphical, inference, prediction)
  - Clinical, manufacturing & commercial informatics

- **Challenge & opportunity:** interdisciplinary company with need for cross-domain expertise *and* platforms that integrate and connect



ACGTTGCATCAGTCAGTCAC

deoxyadenosine monophosphate

d(A)p

p        d        A

# Blending KNIME with Common Drug Discovery Tools



Collectively, this developing system provides curated data and an electronic audit trail that supports IP and regulatory filings *and* quick and easy access to FAIRified data for deeper analytics.
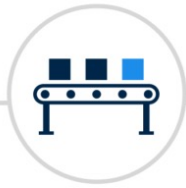
# Industry 4.0 & Pharma 4.0 Framework

# Our KNIME Setup

## KNIME Server

- AWS-hosted
- Test and Prod environments
- Cron job, microservice and WebPortal workflows

## KNIME WebPortal
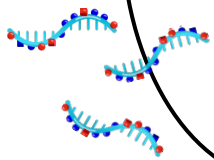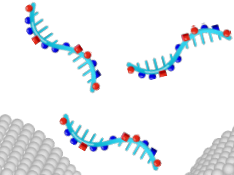
- Heavy internal use by scientists/drug discovery programs
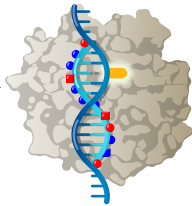
## Integrations

- R
- Python

WAVE
LIFE SCIENCES

# Unlocking RNA editing with PRISM platform to develop AIMers: A-to-I editing oligonucleotides
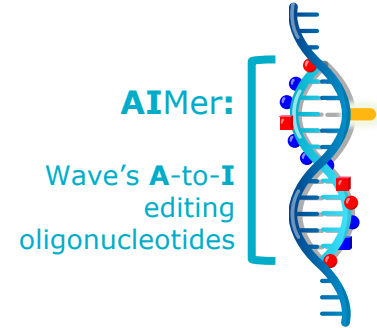
**Free-uptake of chemically modified oligonucleotides**

- First publication (1995) using oligonucleotide to edit RNA with endogenous ADAR[1]

- Wave goal: Expand toolkit to include editing by unlocking ADAR with PRISM oligonucleotides

- ✓ Learnings from biological concepts

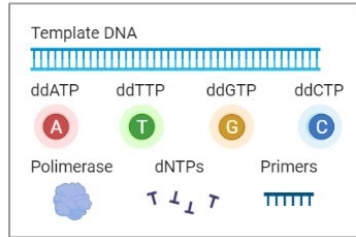- ✓ Applied to ASO structural concepts

- ✓ Applied PRISM chemistry

**AI**Mer:

Wave's **A**-to-**I** editing oligonucleotides

**Endogenous enzymes**

*ADAR*
*RNase H*
*AGO2*
*Spliceosome*

**ADAR enzymes**

- Catalyze conversion of A-to-I (G) in double-stranded RNA substrates

- A-to-I (G) edits are one of the most common post-transcriptional modifications

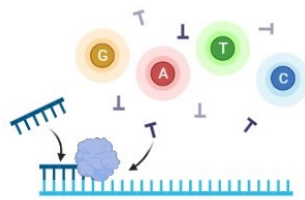- ADAR1 is ubiquitously expressed across tissues, including liver and CNS

[1]Woolf et al., PNAS Vol. 92, pp. 8298-8302, 1995

'WAVE®
LLIFE SCIENCES

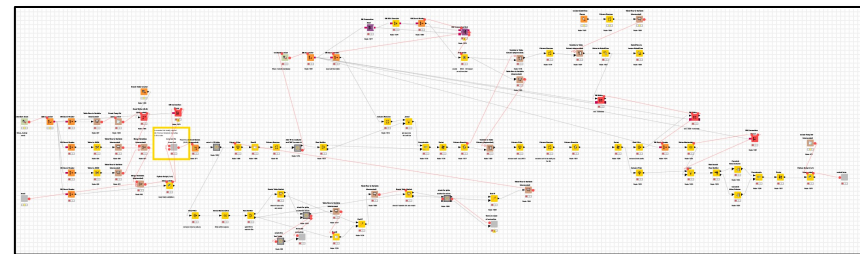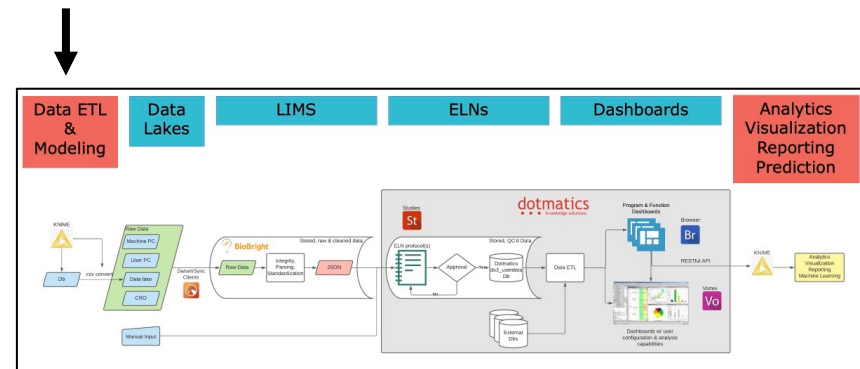# Base Editing Pipeline Analyzes Sanger Sequencing Data



- Sanger sequencing technology has been around in various forms since 1977

- Modern applications produce a sequence chromatogram "ab1" file

- The colors of the chromatogram are translated into the "AGCT" letter code

- We infer the %base editing from the AUCs of the traces

# Sanger Base-Editing Pipeline

The **Sanger pipeline** for **RNA A->I editing measurement** is comprised of three KNIME workflows:
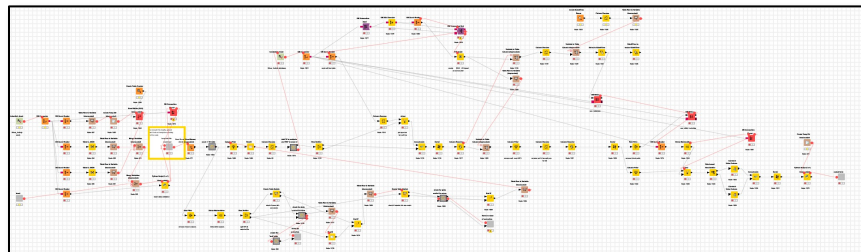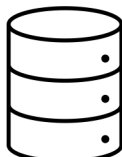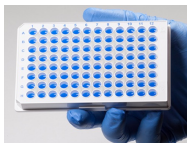
- **Metadata Input and Submission File Generation**
  - 96-well plate format input file to store metadata
  - Exportation of vendor Sample Submission Form

- **Automated Data QC and Processing**
  - CRON job
  - QC of Ab1 files
  - Base representation at all transcript locations along AIMer calculated via editR

- **Data Retrieval**
  - Export edit site percent editing results
  - Download chromatograms

# Metadata Registration and Submission File Generation (1)



User input



Standardized experimental metadata with built-in validation; stored in Db

Python Script

Vendor submission



Automated generation of downloadable vendor-preferred formatted sequencing submission file
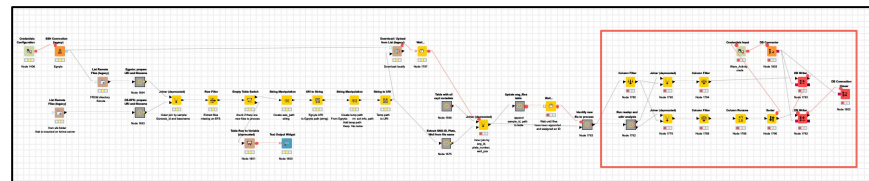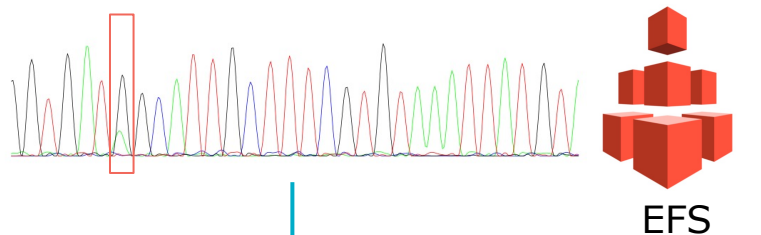
WAVE
LIFE SCIENCES

# Automated ab1 file transfer to Server-mounted EFS (2)



- ab1 files deposited by vendor to Egnyte repository

- CRON job runs on KNIME Server to detect presence of new files in repo

  - Transfer to mounted EFS

  - Store file metadata in database

# Automated Sanger Data QC and Processing (3)



EFS

QC

**R Snippet**

Base Edit Measurement

| S metric | D value | D flag | I file_id |
|---|---|---|---|
| avg_qual | 48.95 | 0 | 26029 |
| read_length | 399 | 0 | 26029 |
| perc_mixed | 0.04 | 0 | 26029 |
| alignment | 1 | 0 | 26029 |
| avg_qual | 52.636 | 0 | 26030 |
| read_length | 327 | 0 | 26030 |
| perc_mixed | 0.018 | 0 | 26030 |
| alignment | 1 | 0 | 26030 |
| avg_qual | 49.28 | 0 | 26031 |

| I file_id | I aso_index | D edit_site | D A_perc | D C_perc | D G_perc | D T_perc |
|---|---|---|---|---|---|---|
| 26027 | 30 | 0 | 50.373 | 1.866 | 47.761 | 0 |
| 26027 | 29 | 0 | 2.101 | 92.437 | 0.42 | 5.042 |
| 26027 | 28 | 0 | 1.676 | 32.682 | 7.263 | 58.38 |
| 26027 | 27 | 0 | 15.134 | 1.484 | 78.932 | 4.451 |
| 26027 | 26 | 0 | 9.091 | 83.636 | 3.182 | 4.091 |
| 26027 | 25 | 0 | 7.725 | 84.549 | 0 | 7.725 |
| 26027 | 24 | 1 | 63.959 | 3.553 | 26.396 | 6.091 |
| 26027 | 23 | 0 | 59.162 | 0 | 33.508 | 7.33 |
| 26027 | 22 | 0 | 56.693 | 4.331 | 34.646 | 4.331 |

- ab1 files QC'ed based on
  - Quality metrics
  - Compound alignment to sequence
- %base editing calculated across entire compound

WAVE
LIFE SCIENCES

# Data Retrieval (4)



- **WebPortal** downloadable Excel files containing experimental metadata, base editing statistics and chromatogram plots

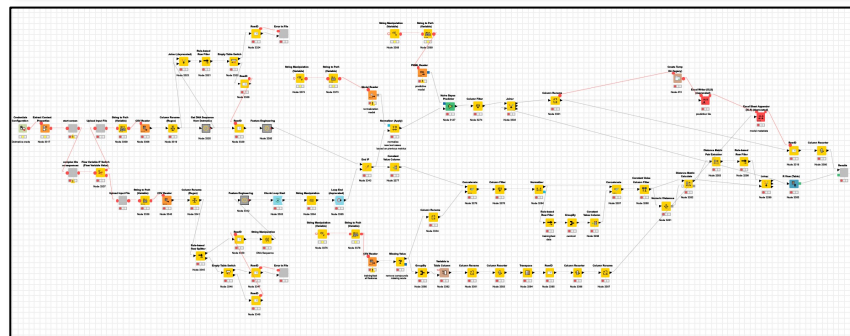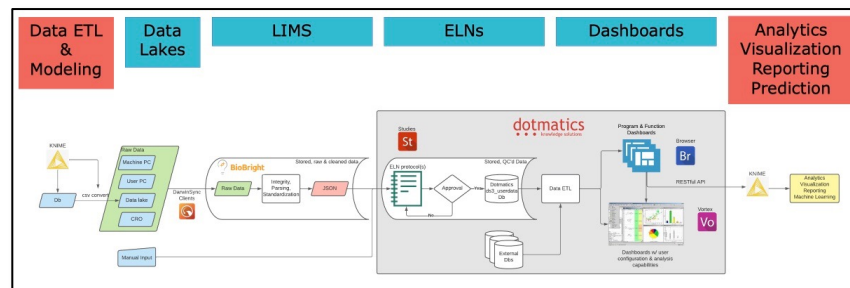- **Dotmatics** Browser-based dashboard for program teams

# Measuring (& Predicting) Tolerability in Mice

We **monitor behavior** of mice in response to chemical compounds repeatedly over several weeks:
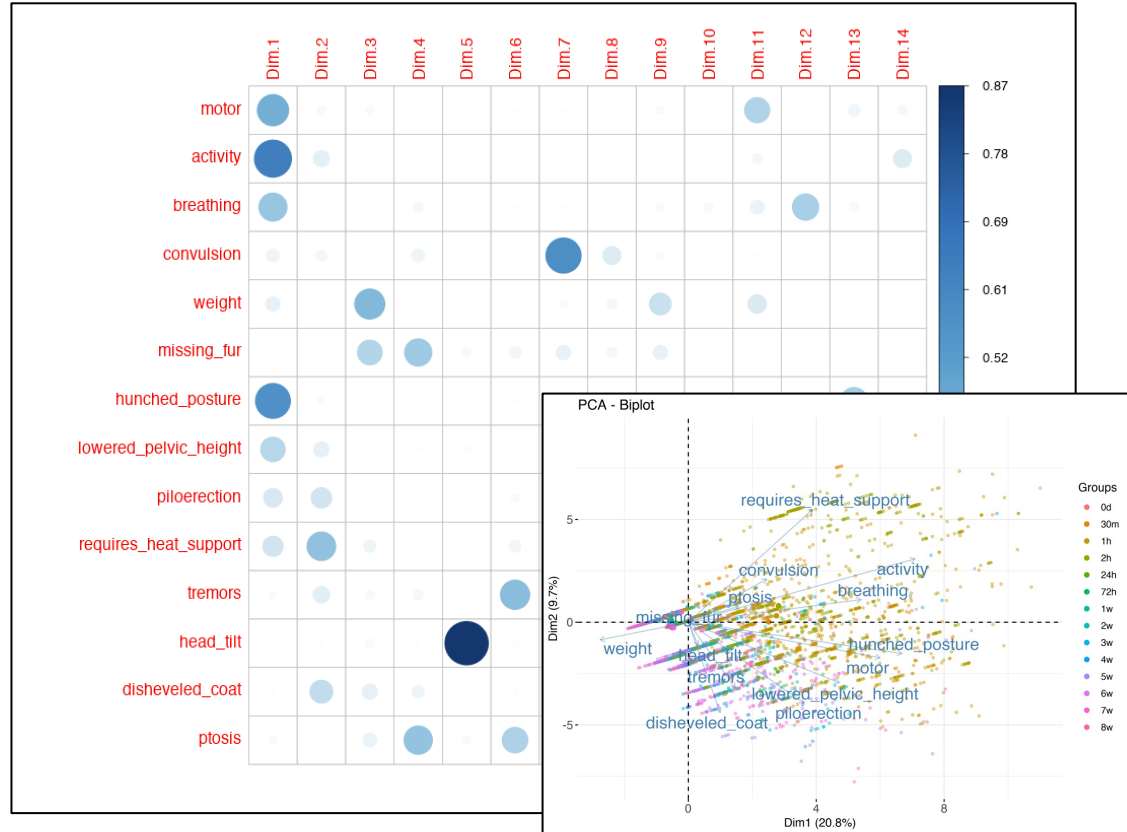
- Activity

- Motor changes

- Breathing

- Physical characteristics

- Body weight

We use these data to build **predictive models** that help to inform our medicinal chemistry design and make molecules safer.

# Behavioral Features Using Principal Components

- We used **PCA** to analyze data from mice **treated with molecules**

- PCA reduces a large number of complex variables and responses (behavior, body weight, etc) to a **smaller number of patterns**

- We ask **how different treatments influence these patterns**, and their relationships to chemical features of molecules
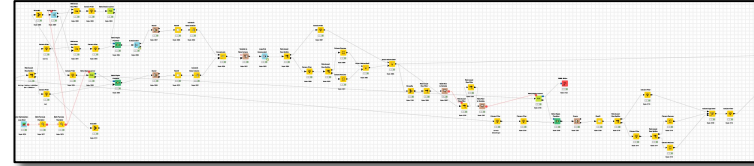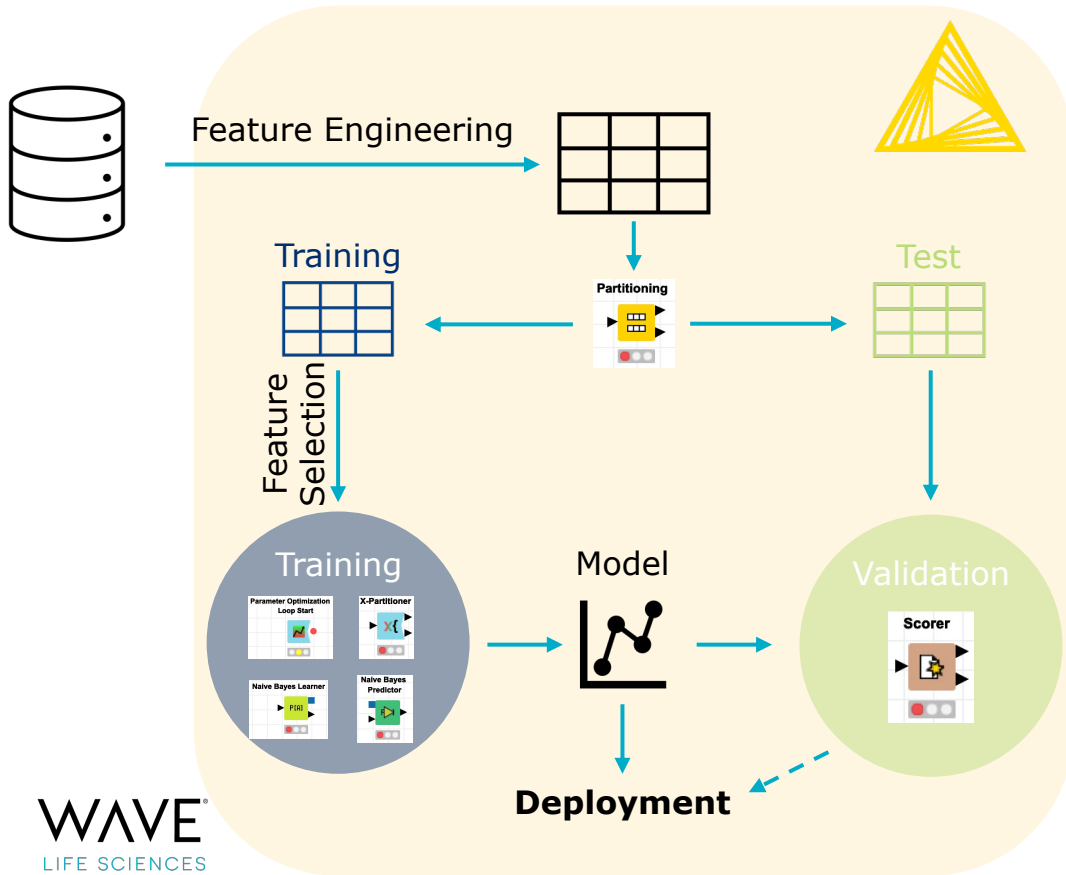


PCA - Biplot

# KNIME Server Webportal GUI for *in vivo* Wild-type Tolerability Prediction

The **WT Tolerability Predictive pipeline** uses KNIME for:

- **Predictive Model Building and Validation**
  - Feature Engineering and Selection
  - Model training and validation

- **User-operated WebPortal tool for predictions**
  - Fast calculation of predictions for new molecules
  - Distance measurement between new cases and existing model dataset

# Predictive Model Building and Validation



- KNIME utilized for data cleaning, feature engineering/selection, and model building and evaluation

- Naïve Bayes Model built and validated for *in vivo* **tolerability prediction**

- 83% accuracy

# Web Tool for Tolerability Prediction of New Compounds



- User submits list of novel molecules

- Model returns tolerability score predictions

- Evaluate distance of predictions from historical tolerated molecule projections

# Summary & Areas for Further Development

- Internal use of component building & sharing

- Expand our microservices framework and RESTful access to data

- Continue to refine our continuous integration/deployment framework

- Expand the reach of KNIME to non-data scientists through internal training programs

- Continue to develop along an Industry/Pharma 4.0 trajectory

WAVE®
LIFE SCIENCES

Thank you!

WAVE®

LIFE SCIENCES