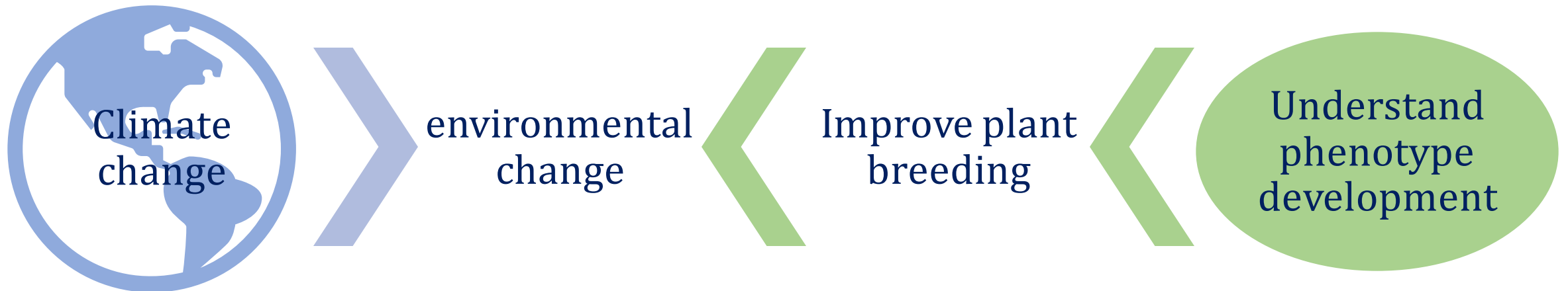


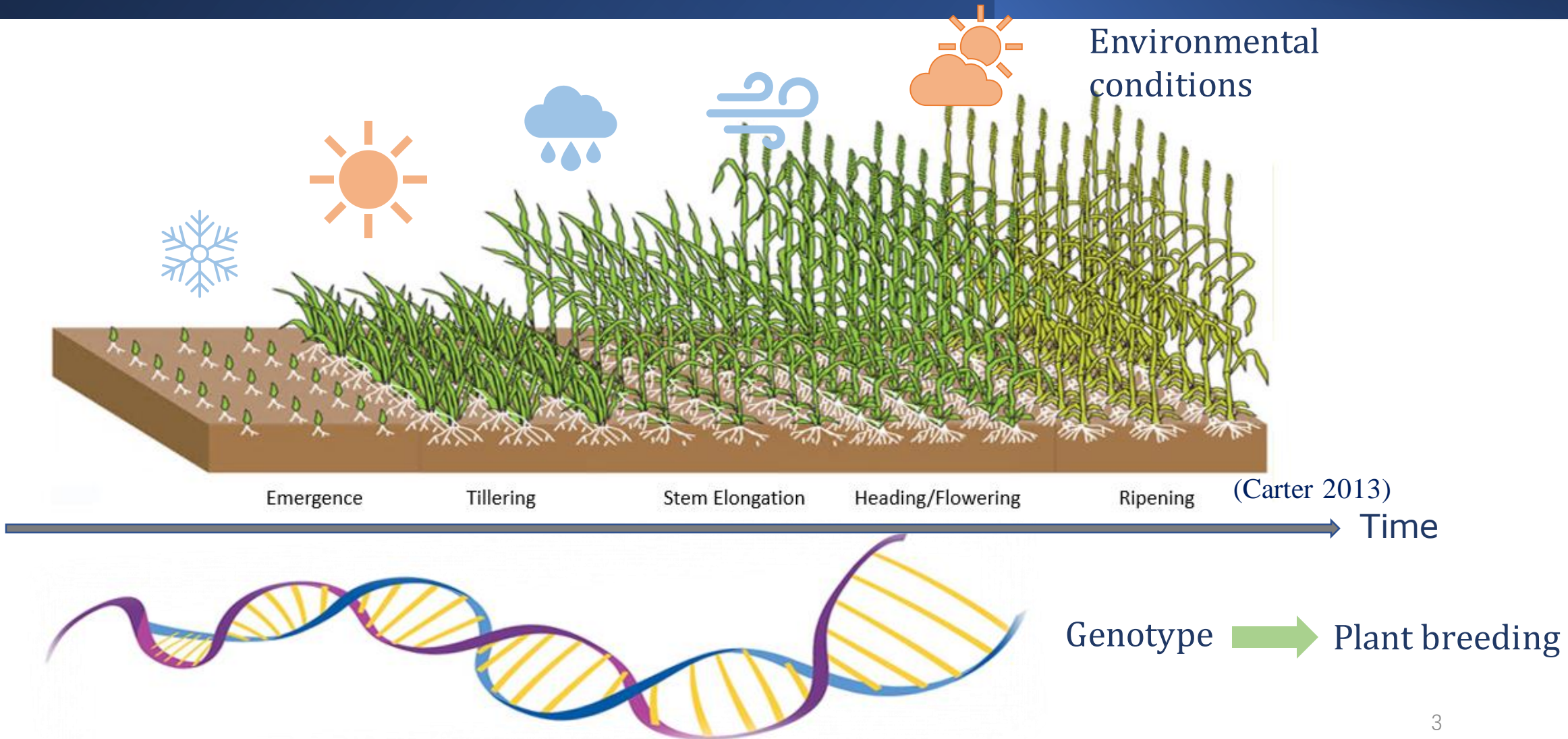
Machine learning for candidate crop growth model classification

Yingjie Shao

Understanding phenotype development will help plant breeding



Phenotype = f(Genotype, Environment)

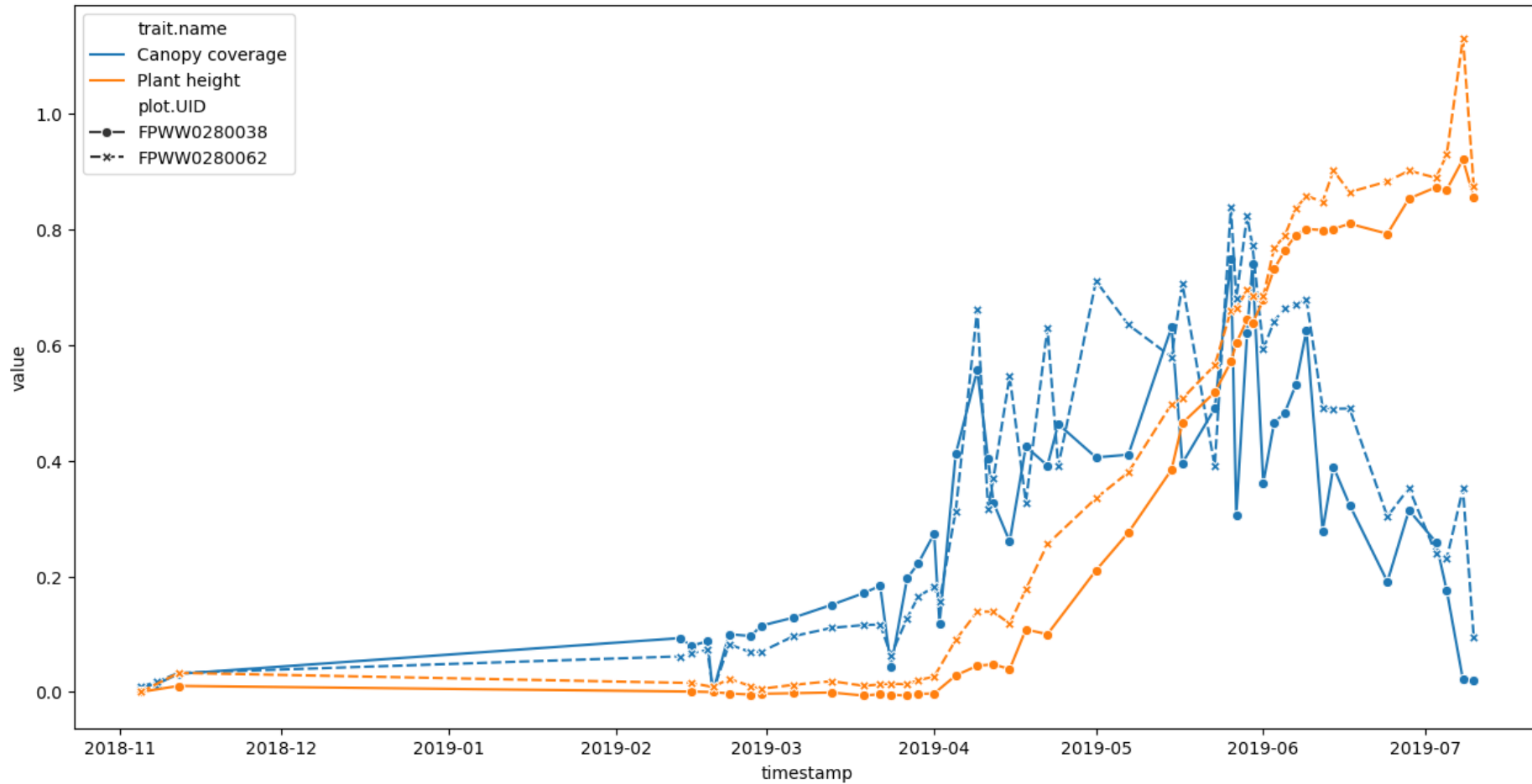


High-throughput phenotyping (HTP) Data

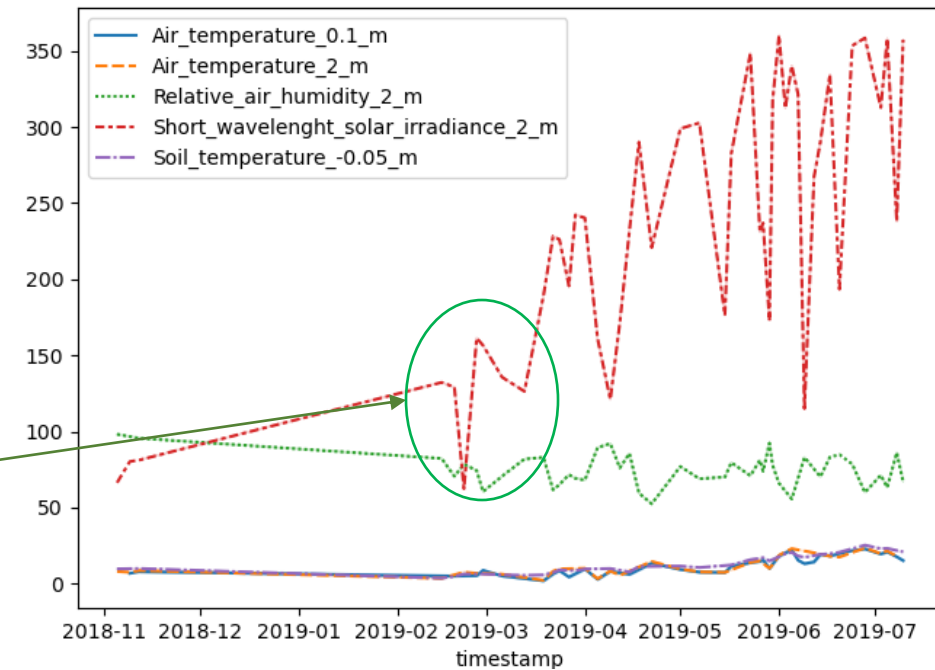
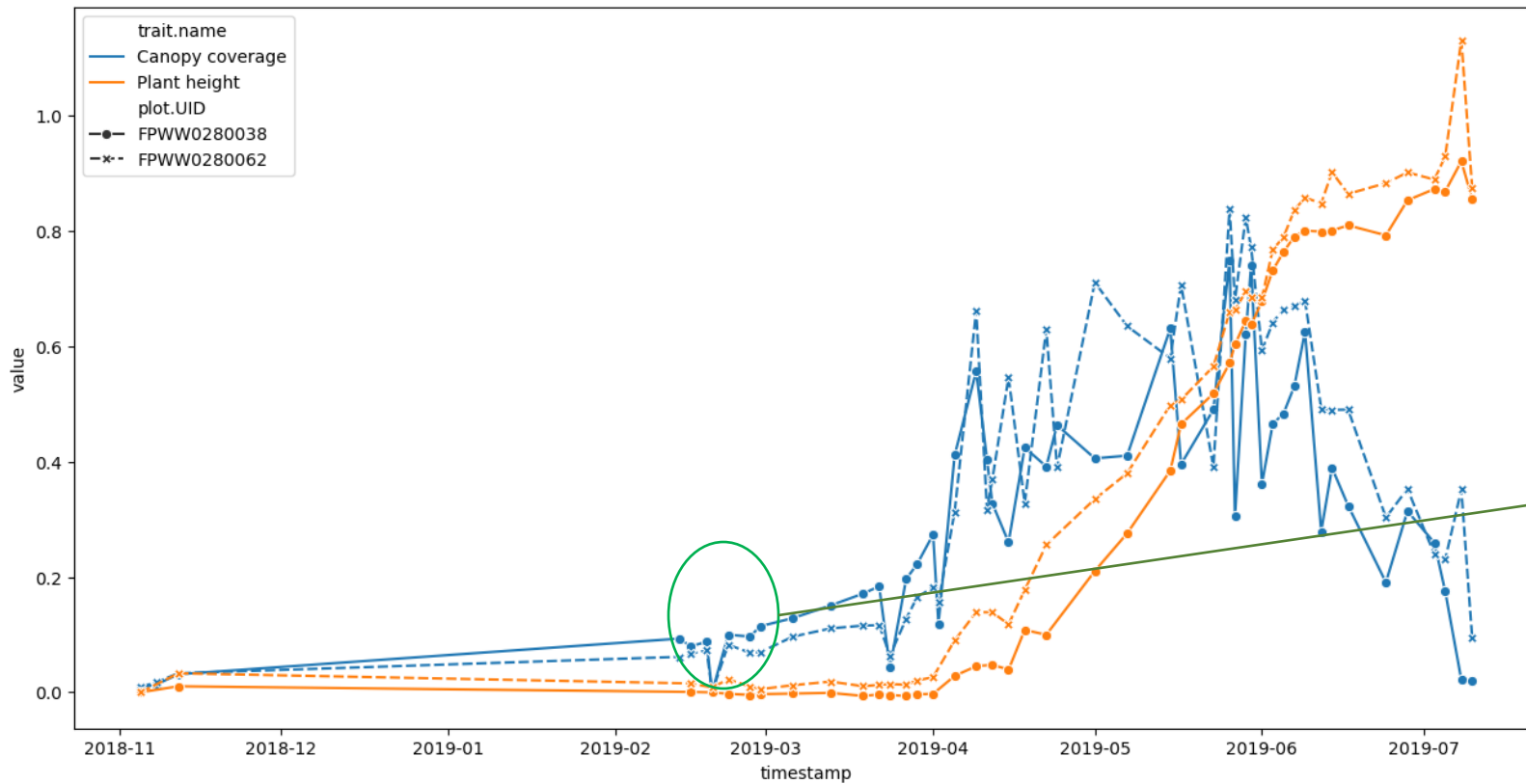
- Non-destructive phenotyping
- Produces multiple traits time series



Time Series Information in Plant Growth



Why Machine learning (ML)?



Extract patterns from time series (crop growth curves), map to non-linear environment effects

Aim

To link phenotype development with the changing environment...

In my MSc thesis, we test how ML works in a simpler case:

Can we use machine learning models to classify the crop growth curve into classes with different environmental limitations?

ODE model for dynamic crop modelling

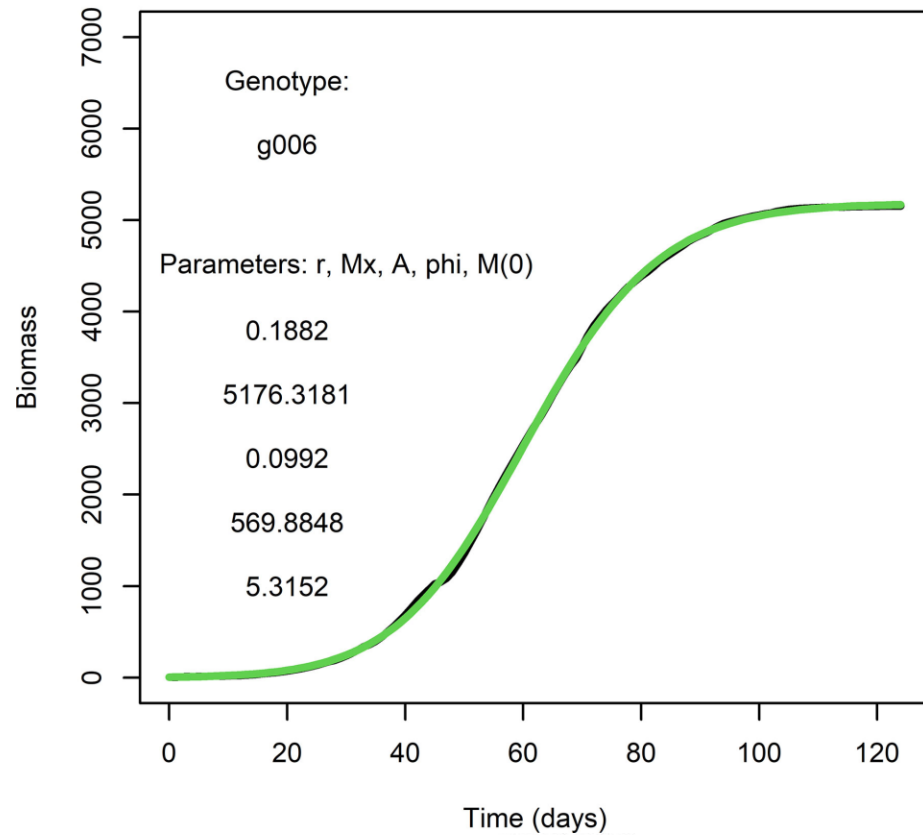
- To generate simulated data...
- Model crop growth as derivatives with respect to time
- Describe dynamic non-linear system

Four types of growth model simulations

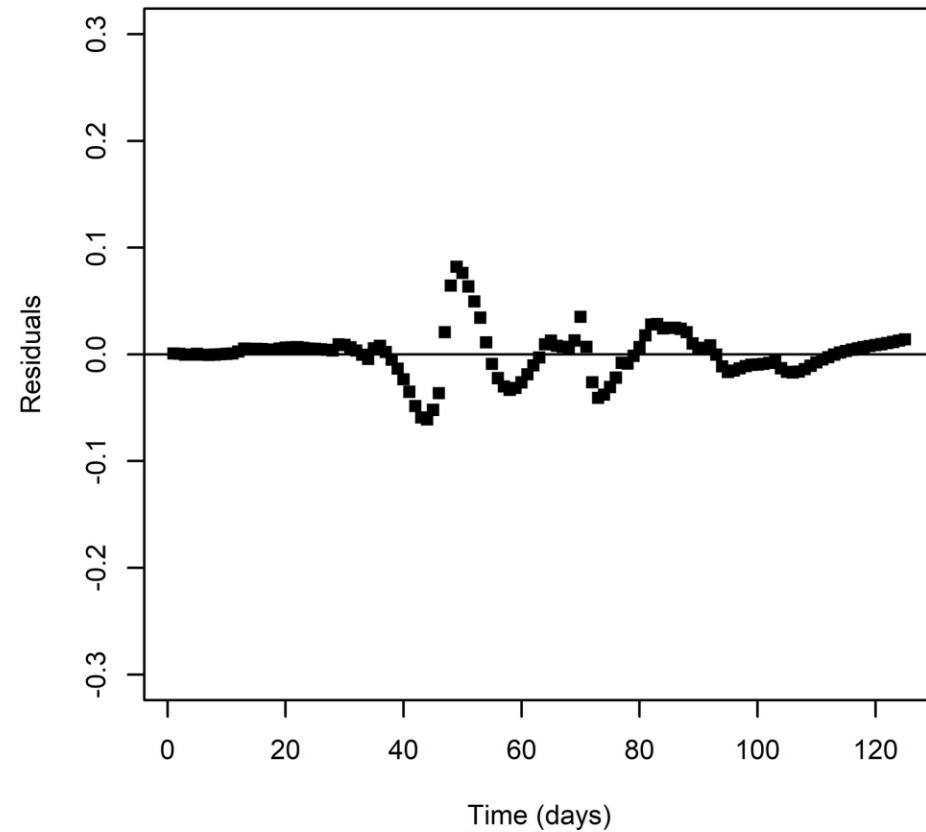
Growth Model*	Formula**	Index
Logistic model	$\frac{dM(t)}{dt} = r * M * \left(1 - \frac{M}{M_{max}}\right)$	1
Irradiance model	$\frac{dM(t)}{dt} = \left(r + \left(a * \sin\left(\left(2 * \frac{\pi}{365}\right) * t + \phi\right)\right)\right) * M * \left(1 - \frac{M}{M_{max}}\right)$	2
Temperature model	$r_{adapt} = \left(1 + \left(\left(\exp\left(\frac{T_{AL}}{temp_t + 273} - \frac{T_{AL}}{T_L}\right) + \exp\left(\frac{T_{AH}}{T_H} - \frac{T_{AH}}{temp_t + 273}\right)\right)\right)\right)^{-1}$ $\frac{dM(t)}{dt} = r_{adapt} * r * M * \left(1 - \frac{M}{M_{max}}\right)$	3
Allee model	$\frac{dM(t)}{dt} = r * M * \left(1 - \frac{M}{M_{max}}\right) * \left(\frac{M}{M + Ma}\right)$	4

ODE models describe plant phenotype development

Irradiance model (model #2)



Emerald, 1985

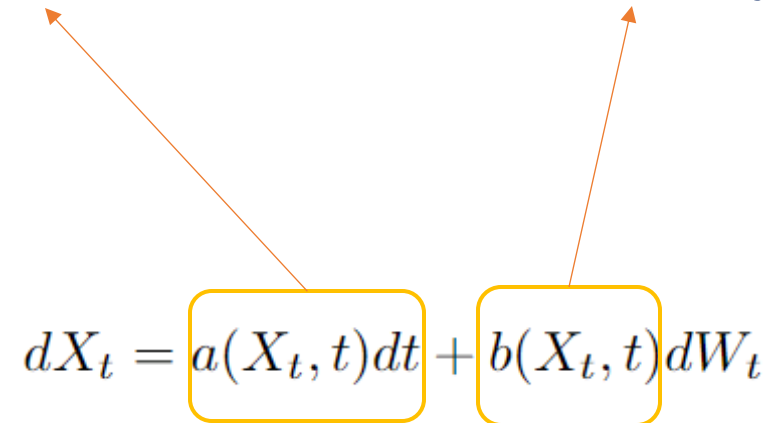


(van Voorn, Boer et al. 2023)

Fit ODE to daily biomass measurements generated from APSIM-Wheat platform

Data Simulation

- Stochastic differential equation (SDE):
 - Uncertainty of plant growth
 - Stochastic arises from environmental factors and measurement techniques
- Different growth models + different noise types

$$dX_t = a(X_t, t)dt + b(X_t, t)dW_t$$
The diagram shows the SDE equation $dX_t = a(X_t, t)dt + b(X_t, t)dW_t$. The terms $a(X_t, t)dt$ and $b(X_t, t)dW_t$ are enclosed in yellow rounded rectangles. Two orange arrows originate from the top-right corner of the first rectangle and the top-left corner of the second rectangle, pointing upwards and outwards towards the text 'Different growth models + different noise types' in the list above.

Four different noise types

Noise type*	Formula
Without noise	$Noise = 0$ Without added noise, it is an Ordinary Differential equation
Time-dependent noise	The highest noise present in the middle of the time series $Noise = 0.2 * \left(\left(2 * \frac{(timelength-t)}{timelength} \right) * \left(1 - \frac{(timelength-t)}{timelength} \right) \right)$
Biomass-dependent noise	The highest noise present when we reach half of the maximum biomass $Noise = 0.2 * \left(\left(2 * (M_{max} - x) / M_{max} \right) * \left(1 - (M_{max} - x) / M_{max} \right) \right)$
Independent noise	$Noise \sim N(0, 0.25)**$

Input Data set

Input X: 4*4 SDE models, 300 samples per SDE model

For LSTM
model

Daily biomass shape = (120,1200,1)

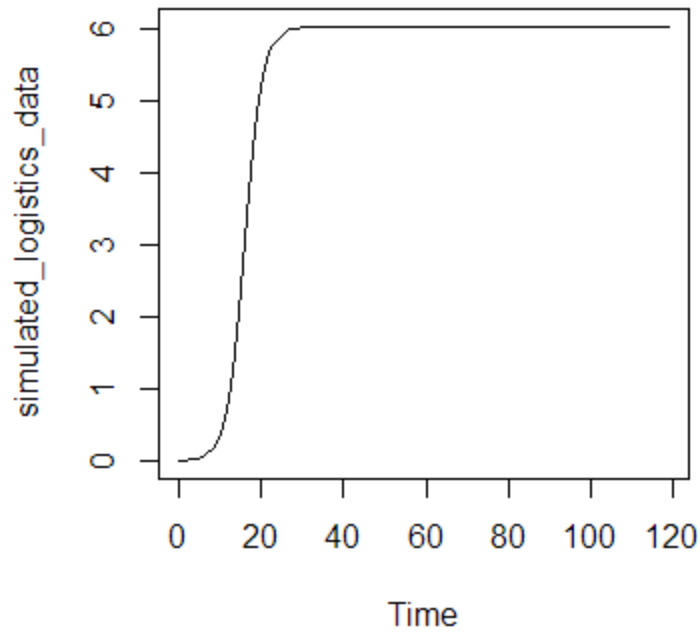
Daily biomass + the derivative of biomass shape = (120,1200,2)

For CNN
model

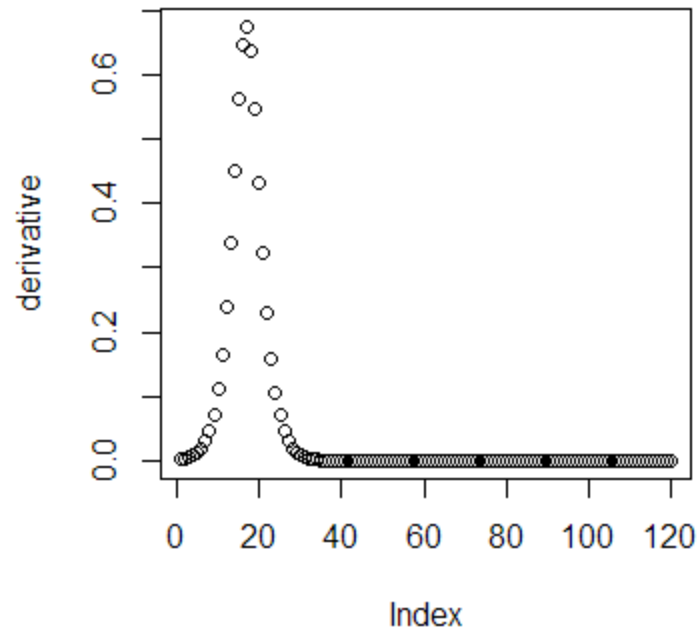
The image of Daily biomass + the derivative of biomass
shape = (1200,2,256,256)

Input Y: growth model types (4 classes labels)

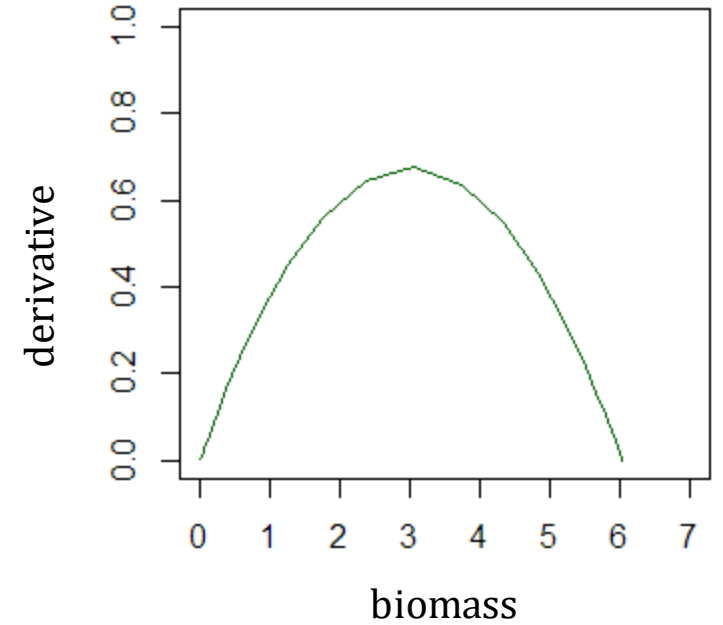
Example of curves without noise



Biomass Growth curve

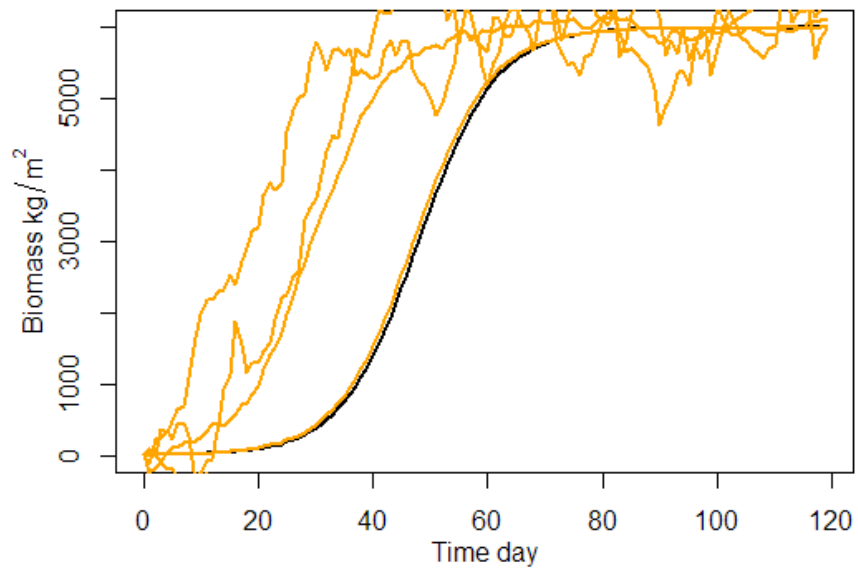


Biomass derivative against days

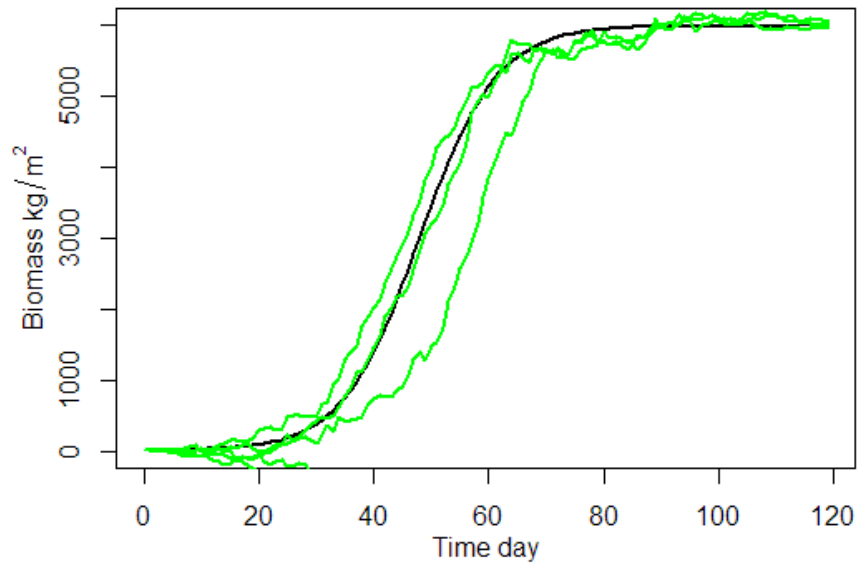


Biomass derivative against biomass

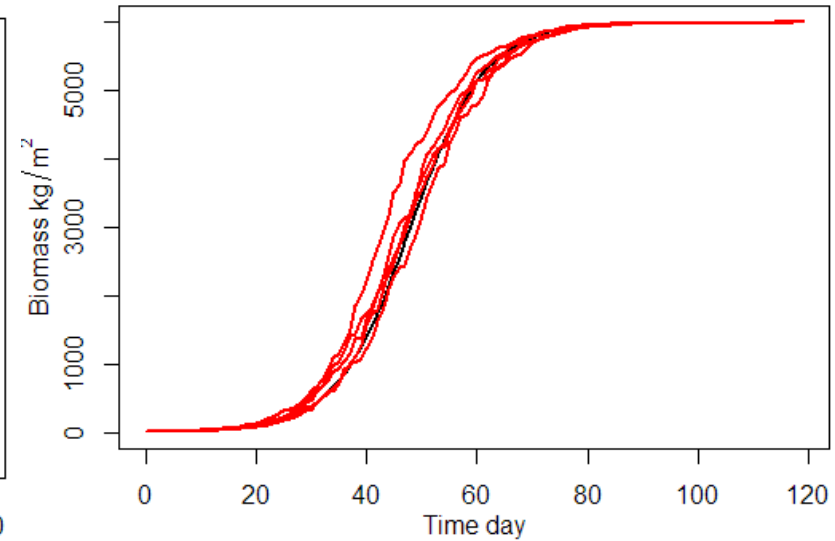
Biomass growth curve for different noise types



Independent noise

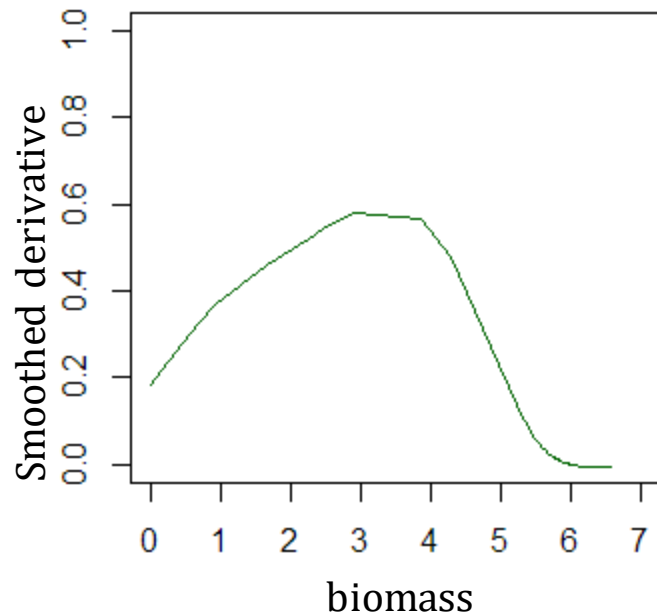


Time-dependent noise

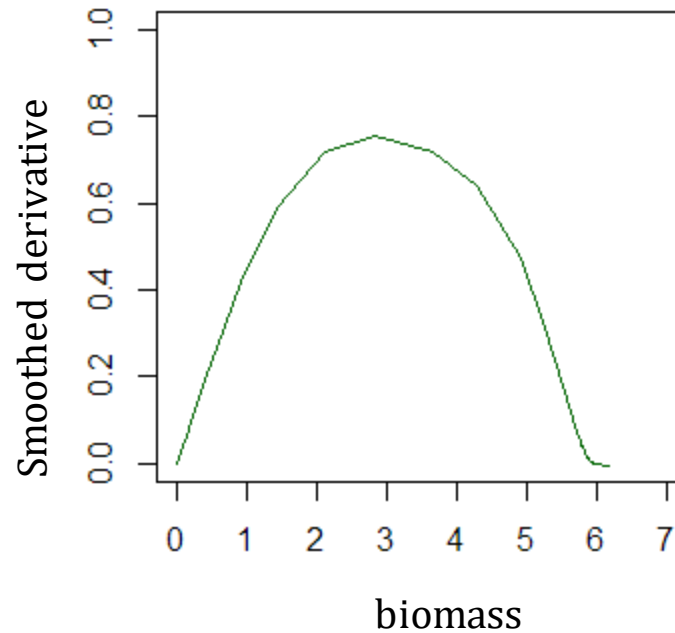


Biomass-dependent noise

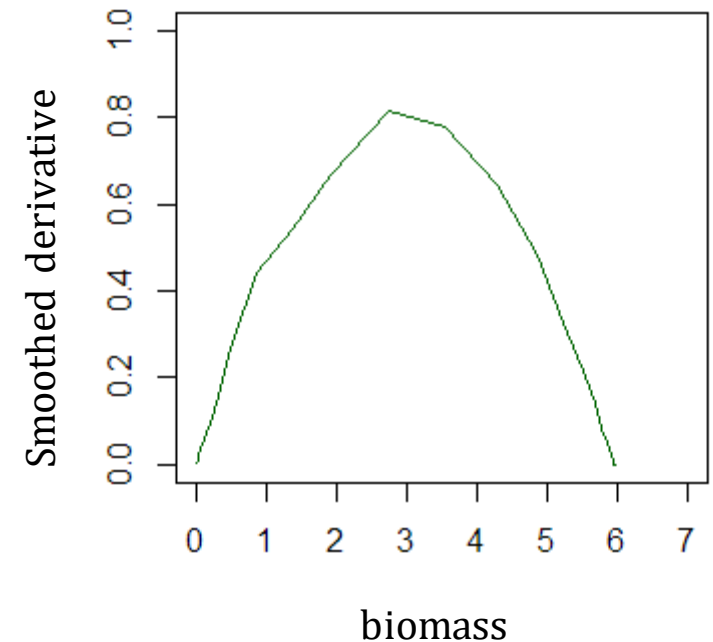
Biomass derivative against biomass for different noise types



Independent noise



Time-dependent noise



Biomass-dependent noise

Input Data preprocessing

Input X:

split based on **noise type** -> models are trained separately for four noise types.

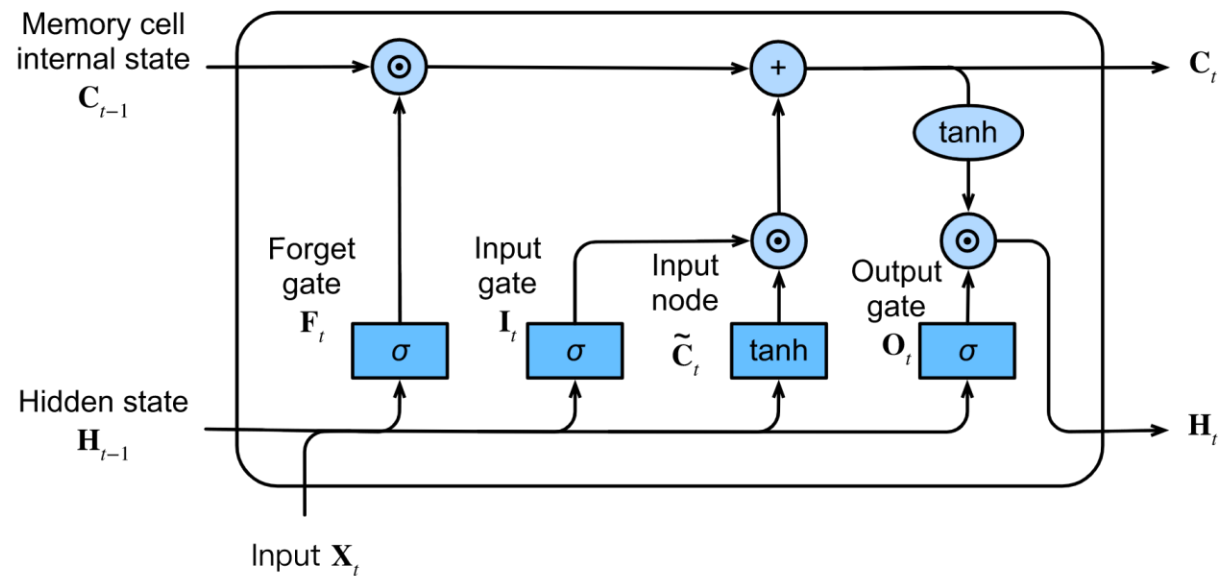
Scaling before feeding into neural networks $z = \frac{x - \mu}{s}$

Input Y: One-hot encoded growth model types

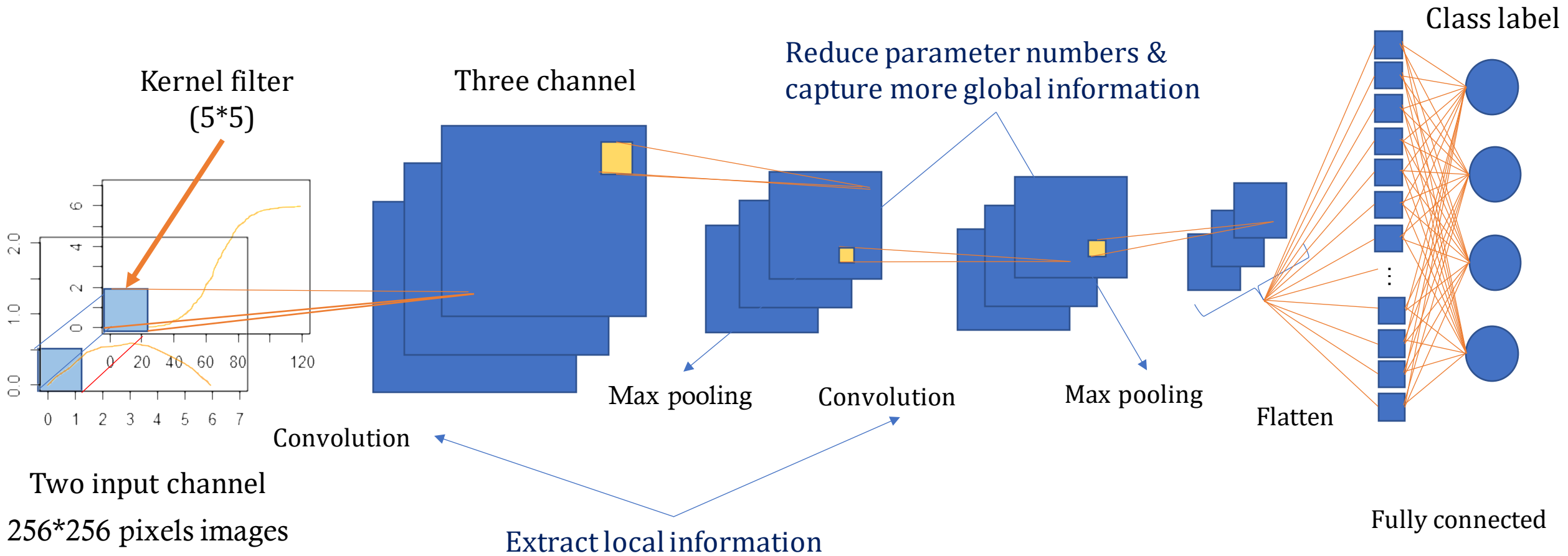
	[1,0,0,0]
	[0,1,0,0]
	[0,0,1,0]
	[0,0,0,1]

LSTM Model

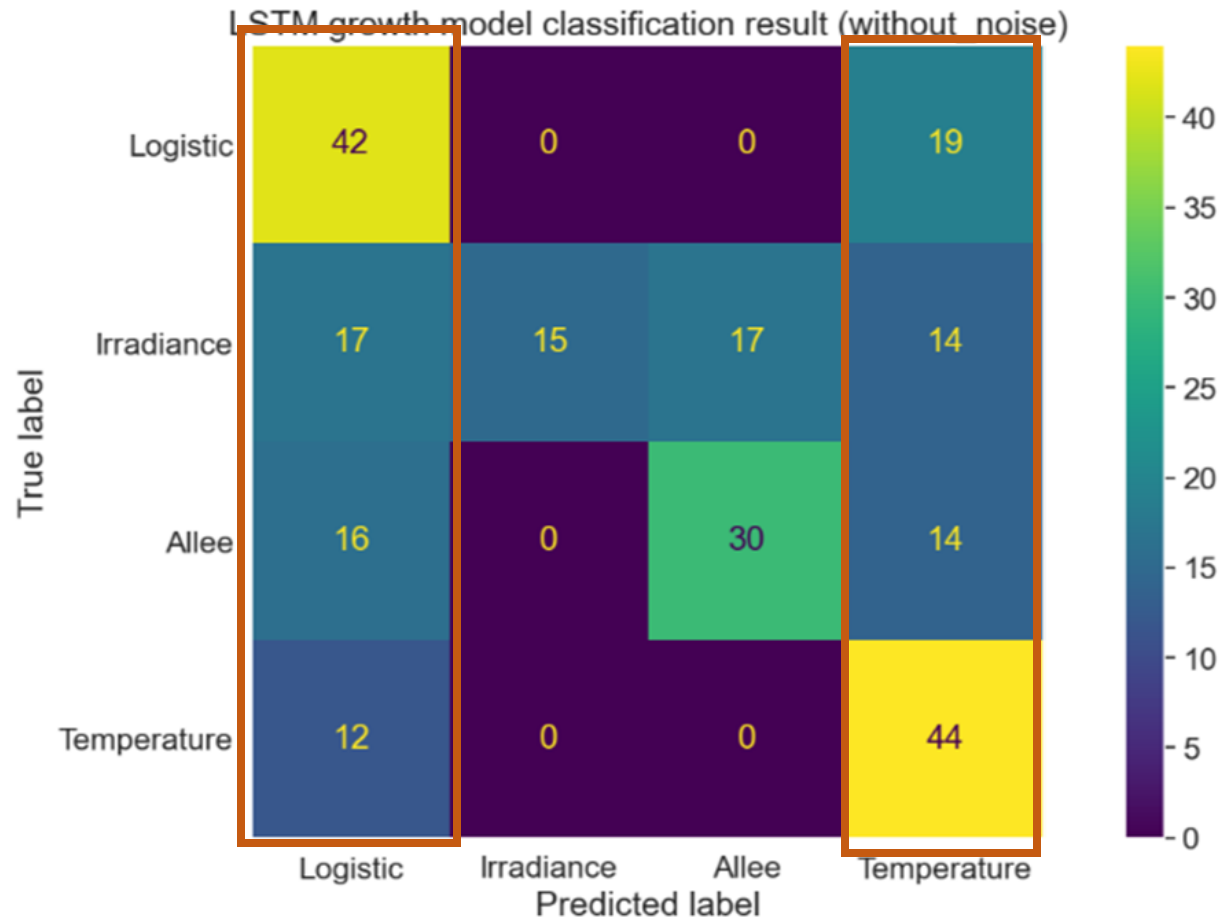
- Two LSTM layer + fully connected layer
- Time series sequences as input



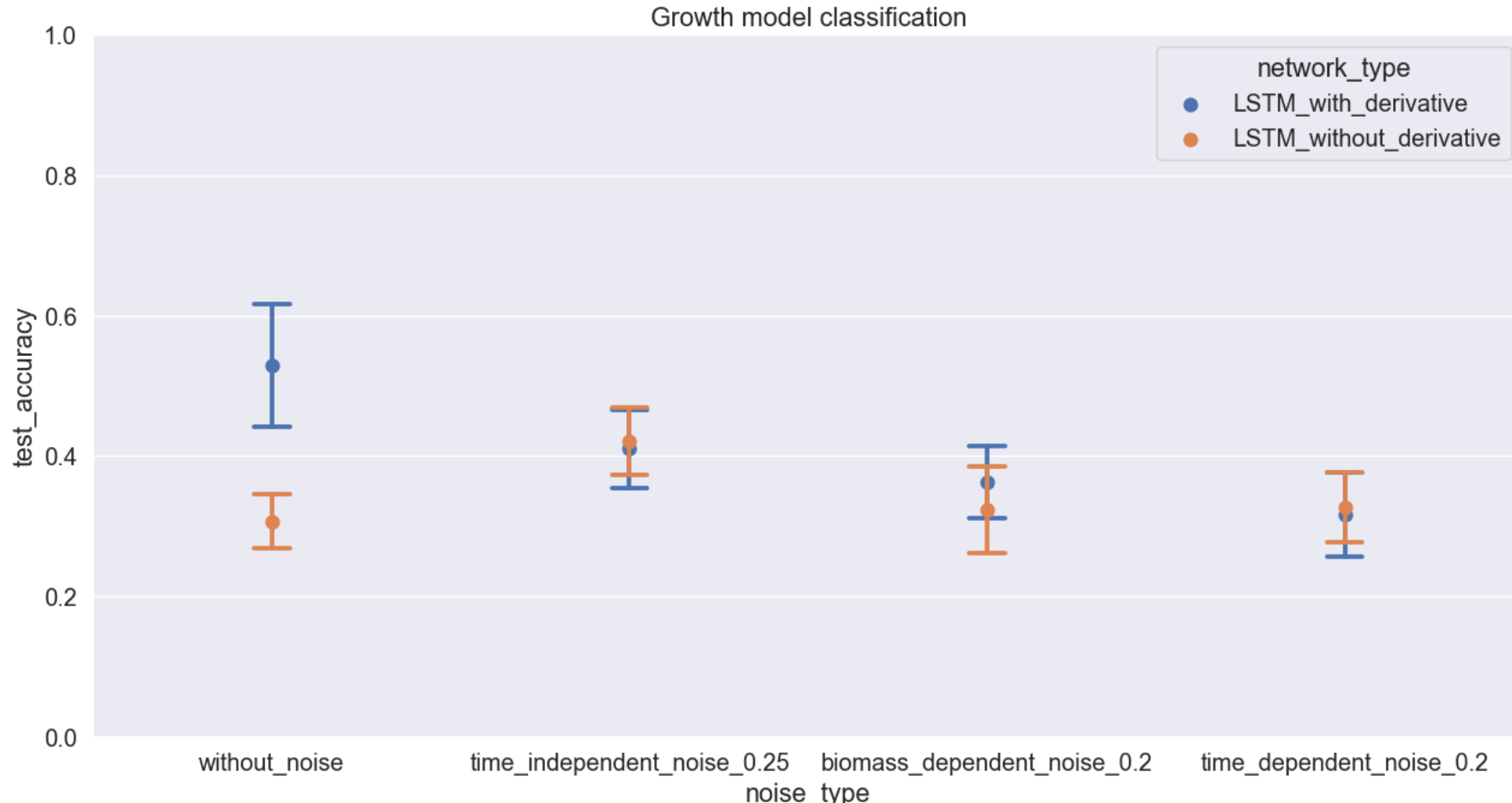
CNN Model



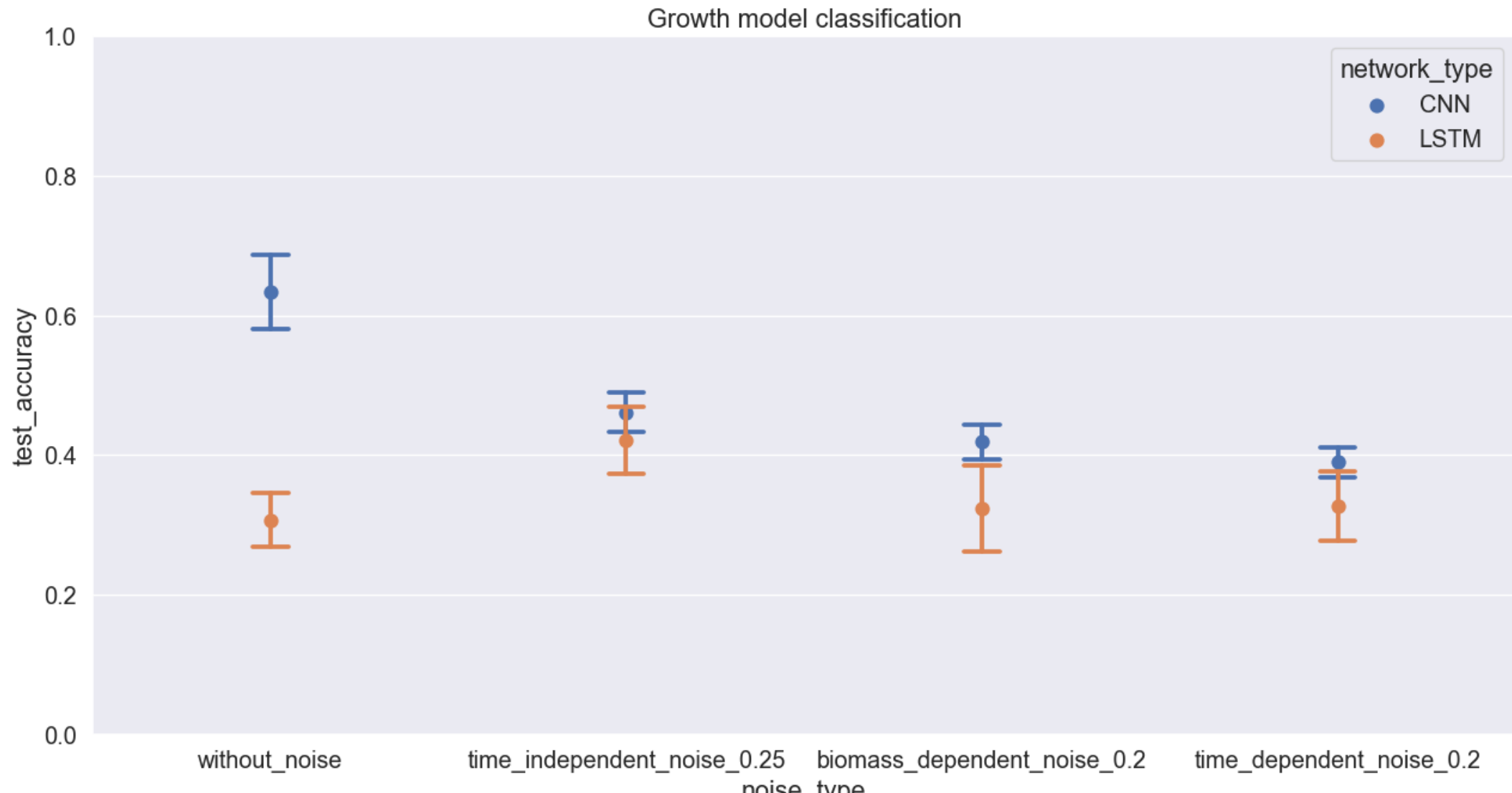
The model tends to classify samples into logistic or temperature class



Added Derivate Information Improves Accuracy



Growth Model Classification Comparison



Conclusion & Discussion



Adding derivative information improves prediction accuracy ->

Both models give relatively high accuracy ->

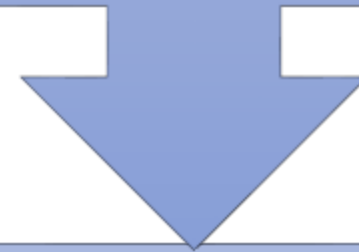


Extracted features from differential equations and combined them with raw data

Combine CNN with LSTM and use previous knowledge to guide hyperparameter selection

Future Plan

Physical knowledge + neural
network



CausalML

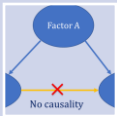
Knowledge gaps



Quantify dynamic genome effects within different environments



Combining process-based and data-driven models in plant studies



Causal discovery and estimation for time series (HTP data)

Objectives

- Build accurate plant growth forecasting models by combining process-based models with neural networks
- Quantify G*E interactions in a dynamic system of plant growth
- Discover causal relationships between phenotypes and environmental factors



Thank you!

Questions?