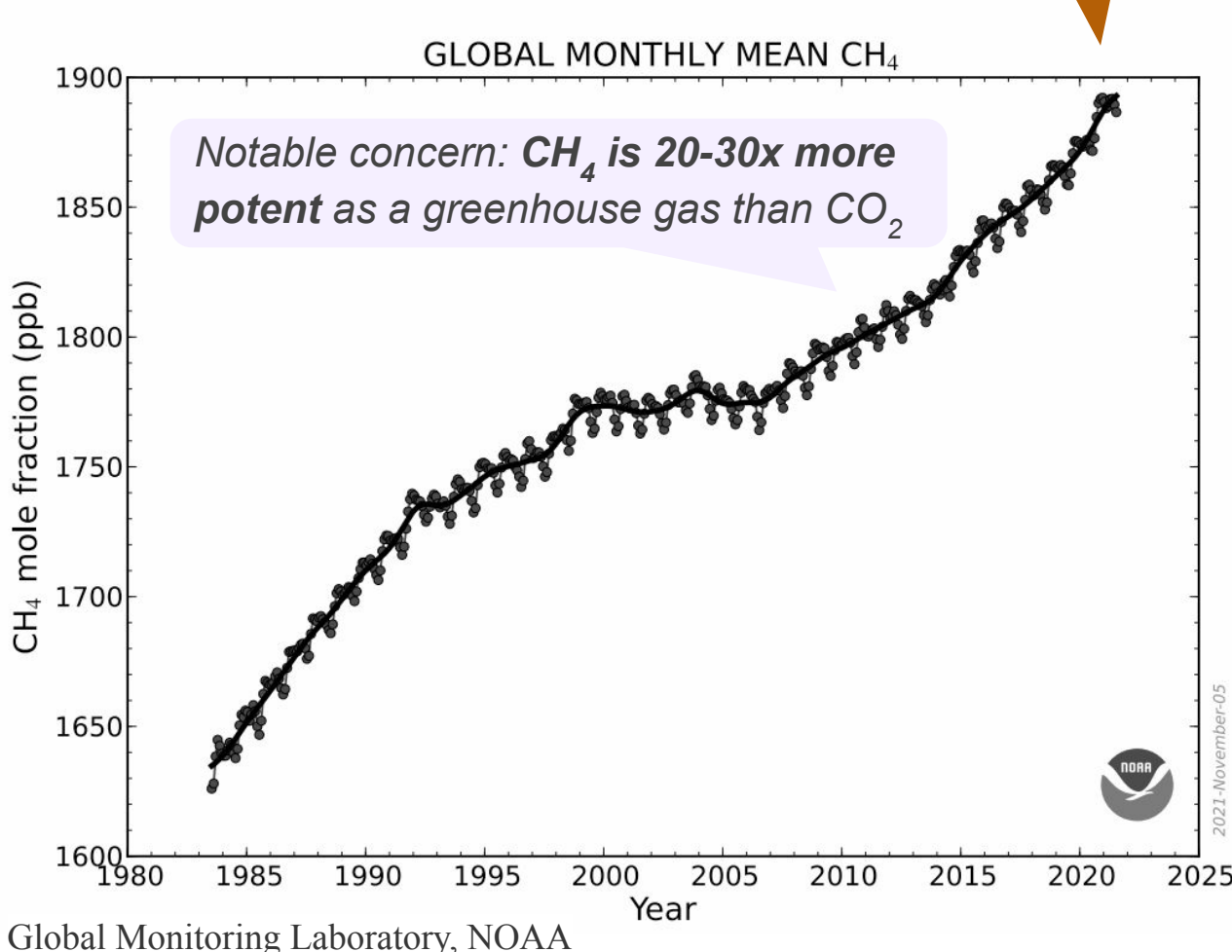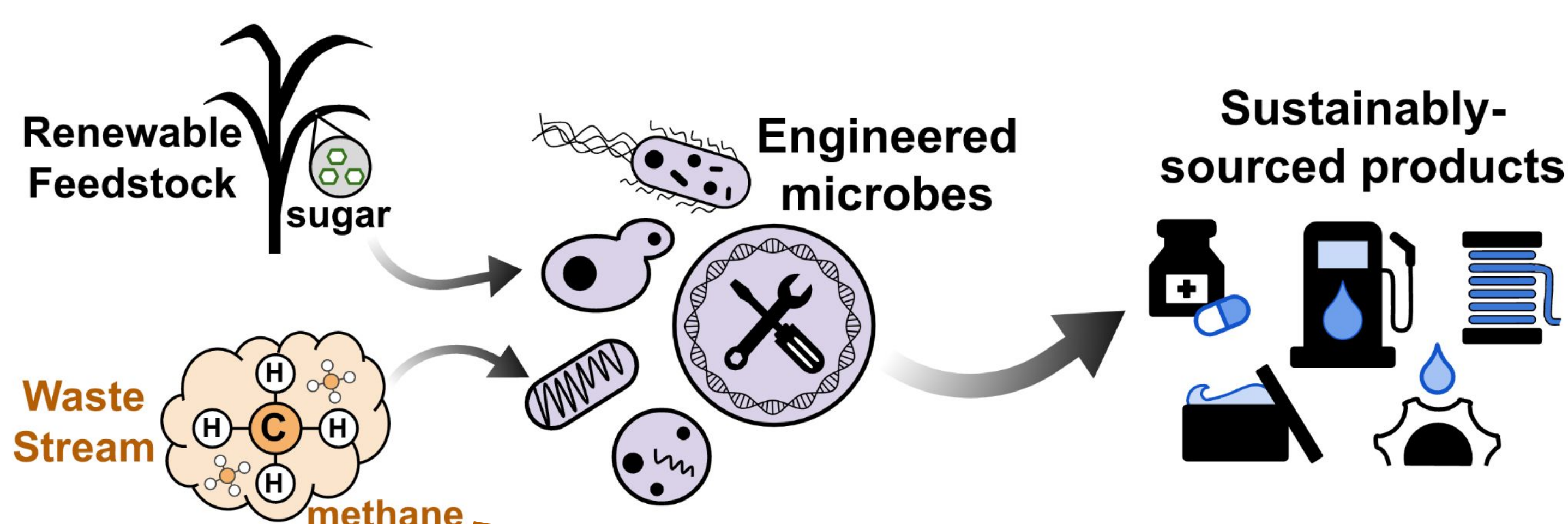# Methane, Microbes, and Machine Learning:
## Engineering biology to combat climate change

**Erin H. Wilson**, Mary E. Lidstrom, David A. C. Beck
ewilson6@cs.washington.edu
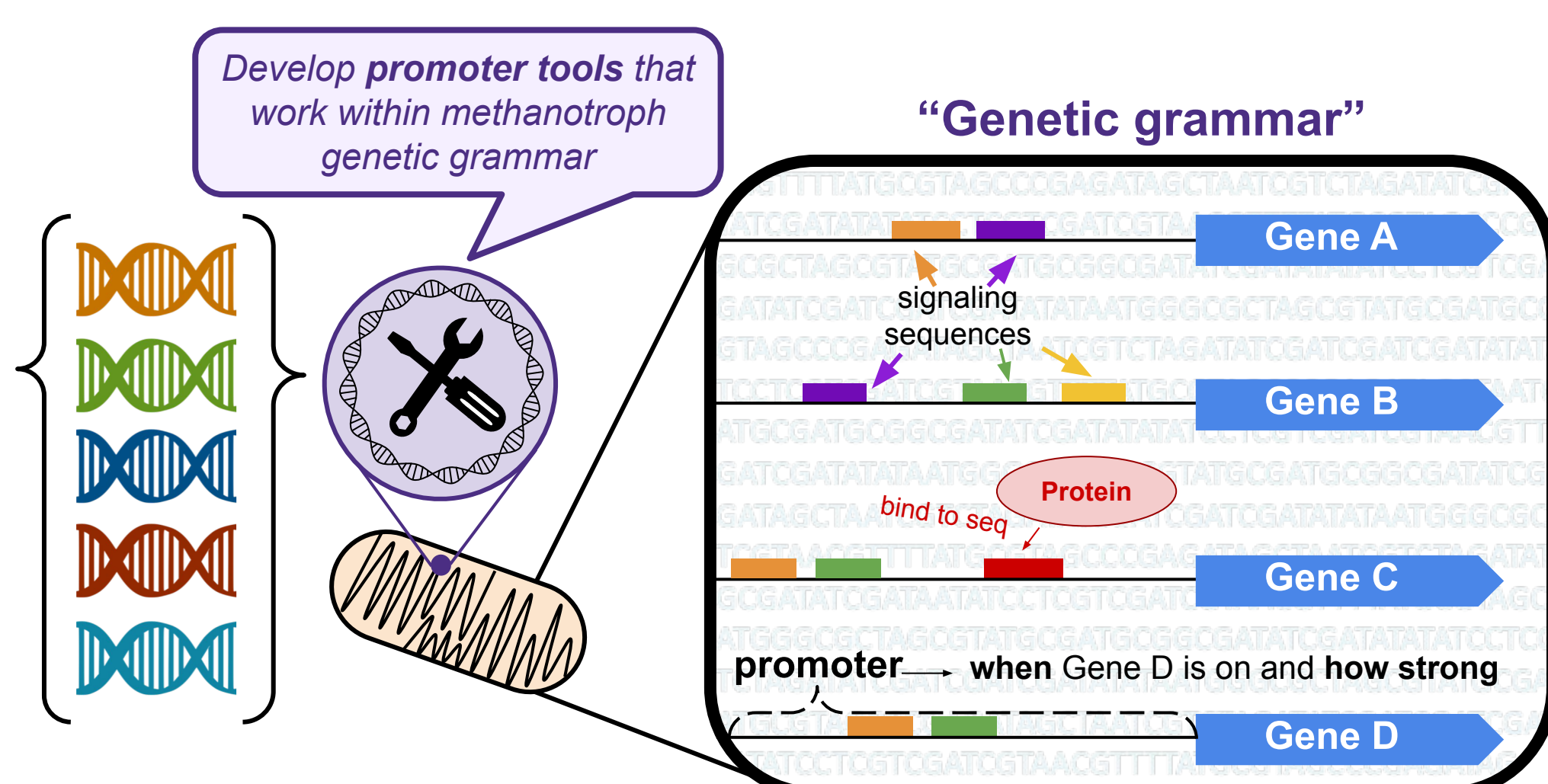
## 1) A promising paradigm for methane mitigation

- Metabolic engineering: a field that aims to **engineer microorganisms into biological factories** that convert renewable feedstocks into valuable biomolecules.
  - → Provides a more **sustainable alternative** to sourcing many materials, especially petroleum-based products
- Much progress with model organisms (baker's yeast and *E. coli*) to produce **malaria medicine, jet fuel,** fragrances



Renewable Feedstock / sugar → Engineered microbes → Sustainably-sourced products

Waste Stream / methane

GLOBAL MONTHLY MEAN CH₄

Notable concern: CH₄ is 20-30x more potent as a greenhouse gas than CO₂

Global Monitoring Laboratory, NOAA

- Methanotrophs - **bacteria** that can **survive on methane** as their sole carbon source - are promising microbial hosts for industrial biomolecule production
- ★ Opportunity to **divert methane waste streams** into **valuable everyday materials**
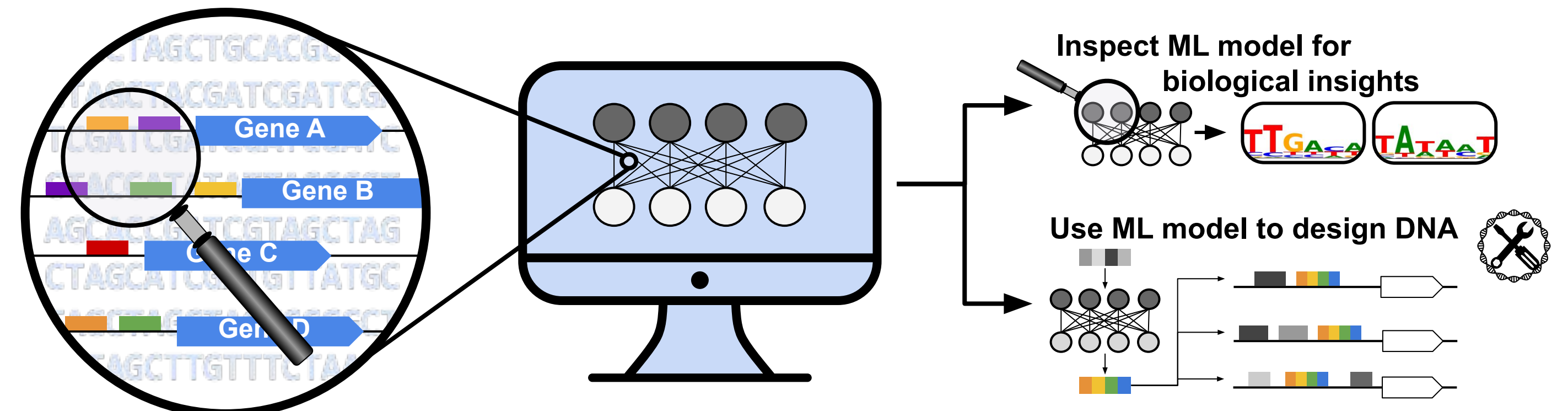
## 2) Regulatory DNA is a complex language to decipher

*Expression of newly installed genes must be carefully controlled in the host microbe*



Develop **promoter tools** that work within methanotroph genetic grammar

"Genetic grammar"

signaling sequences

Protein / bind to seq

promoter — when Gene D is on and how strong
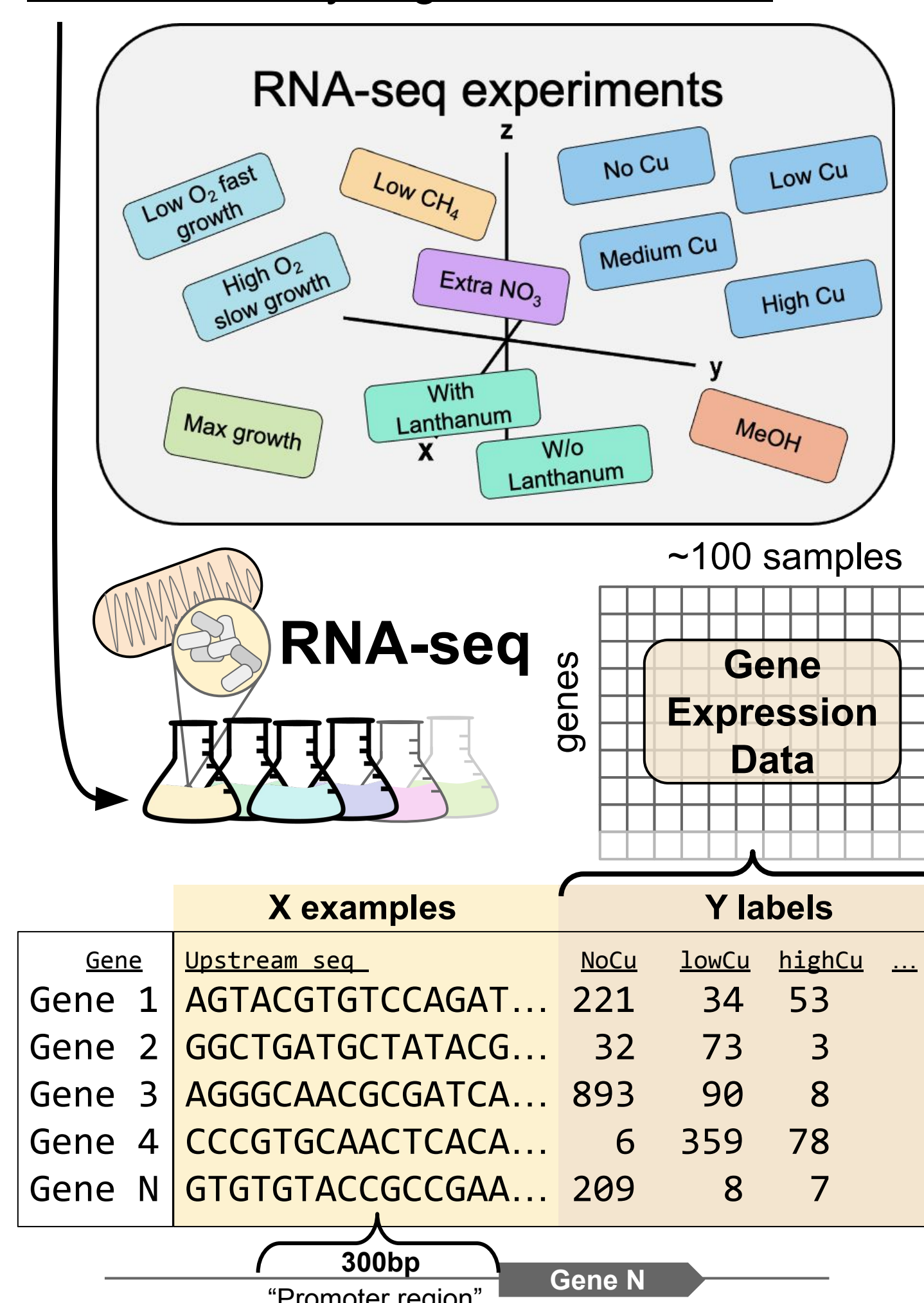
Gene A / Gene B / Gene C / Gene D

- Every microbe has evolved a different **genetic grammar**: a series of **signaling sequences** and logic patterns it uses to control its genes
  - → *Promoter = sequence region containing many signals that influence when genes turn ON or OFF ("expression")*
- We must understand this grammar in order to **efficiently reprogram cells** for biomolecule production
  - → *Research goal: develop methanotroph promoter tools*

## 3) Machine learning to automatically detect patterns in DNA



Inspect ML model for biological insights
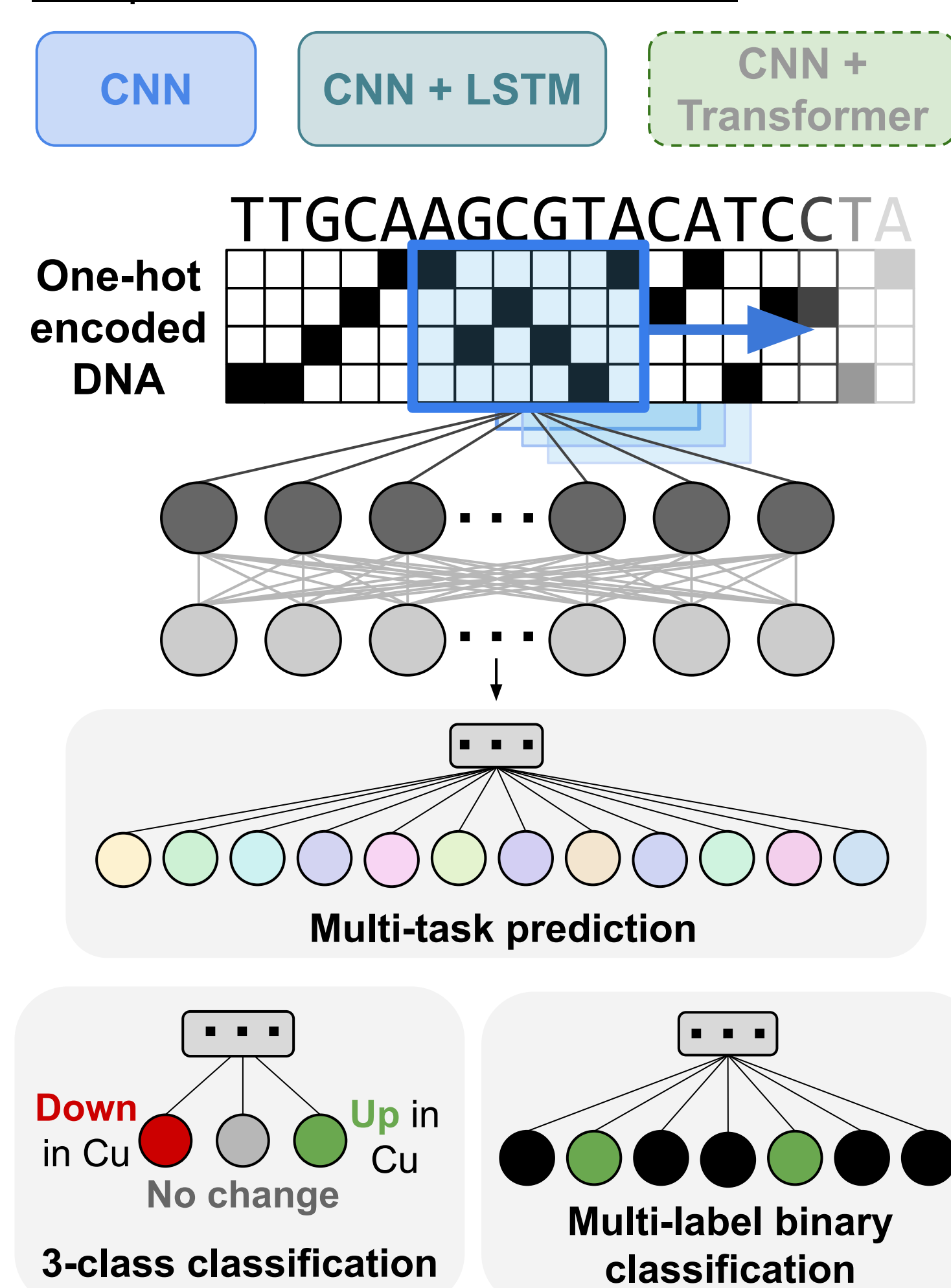
Use ML model to design DNA

- Most DNA sequence signals are still unknown in methanotrophs
- Deep learning approaches can **learn relevant features directly from the data** without explicit encoding
  - → Use deep learning models to **find patterns within methanotroph promoter sequences**
- **Biological insights:** what DNA patterns has the model learned?
- **Novel DNA:** freeze model and use for forward **DNA design**
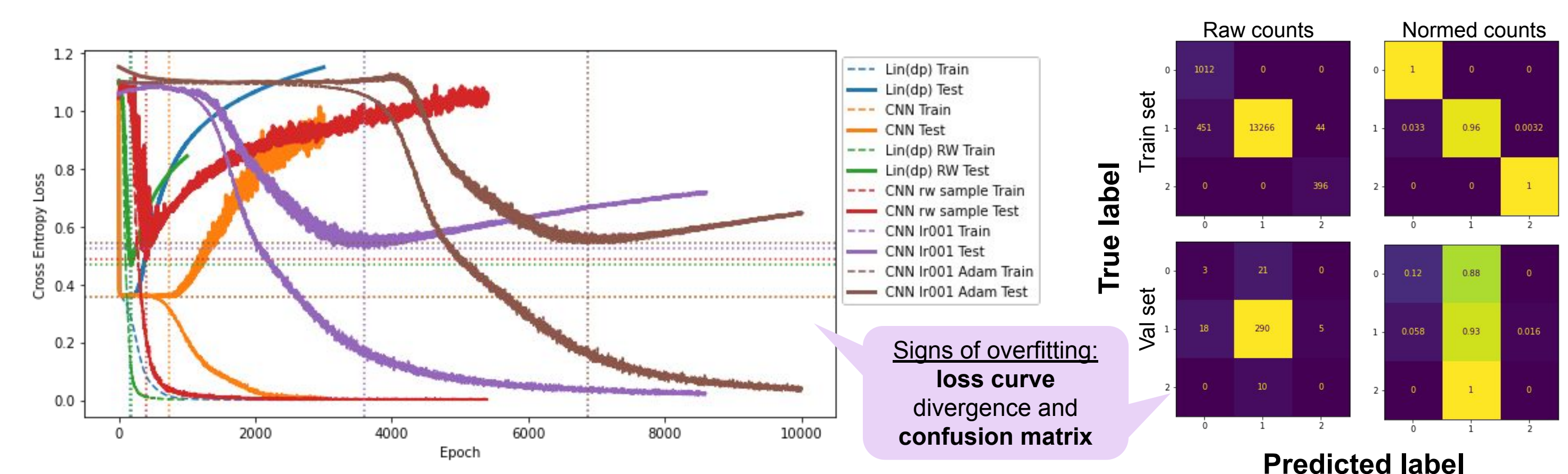
Dataset: variety of growth conditions

RNA-seq experiments

Low O₂ fast growth, High O₂ slow growth, Low CH₄, No Cu, Low Cu, Medium Cu, High Cu, Extra NO₃, Max growth, With Lanthanum, W/o Lanthanum, MeOH

RNA-seq

~100 samples

Gene Expression Data

genes

| Gene | Upstream seq | NoCu | lowCu | highCu | ... |
|------|--------------|------|-------|--------|-----|
| Gene 1 | AGTACGTGTCCAGAT... | 221 | 34 | 53 | |
| Gene 2 | GGCTGATGCTATACG... | 32 | 73 | 3 | |
| Gene 3 | AGGGCAACGCGATCA... | 893 | 90 | 8 | |
| Gene 4 | CCCGTGCAACTCACA... | 6 | 359 | 78 | |
| Gene N | GTGTGTACCGCCGAA... | 209 | 8 | 7 | |

X examples / Y labels

300bp "Promoter region" — Gene N

Sample of ML model architectures

CNN | CNN + LSTM | CNN + Transformer

TTGCAAGCGTACATCCTA

One-hot encoded DNA

Multi-task prediction

**Down** in Cu / **Up** in Cu / No change
**3-class classification**

**Multi-label binary classification**

## 4) Addressing key challenges: overfitting, dataset size, imbalance

- Current models are **overfitting** to the training data, despite initial strategies to address **class imbalance** and **limited data**



Cross Entropy Loss / Epoch

Lin(dp) Train, Lin(dp) Test, CNN Train, CNN Test, Lin(dp) RW Train, Lin(dp) RW Test, CNN rw sample Train, CNN rw sample Test, CNN lr001 Train, CNN lr001 Test, CNN lr001 Adam Train, CNN lr001 Adam Test

Signs of overfitting: loss curve divergence and confusion matrix

Raw counts / Normed counts
Train set / Val set
True label / Predicted label

- Future work: **self-supervised pre-training** on general sequence tasks; fine tune model to methanotroph RNA-seq data

PAUL G. ALLEN SCHOOL OF COMPUTER SCIENCE & ENGINEERING

GRFP / NSF Graduate Research Fellowship Program