

## Using machine learning to model promiscuous activity of thiamine diphosphate-dependent carboligases and side reactions in the *E. coli* metabolome

Tracey Dinh<sup>1\*</sup>([tracey.dinh@u.northwestern.edu](mailto:tracey.dinh@u.northwestern.edu)), Bradley W. Biggs<sup>1</sup>, Matthew T. Robey<sup>2</sup>, Catherine Majors<sup>1</sup>, Lindsay Caesar<sup>2</sup>, Neil L. Kelleher<sup>2,3,4</sup>, Paul M. Thomas<sup>4</sup>, Linda J. Broadbelt<sup>1</sup>, **Keith E.J. Tyo<sup>1</sup>**

<sup>1</sup>Department of Chemical and Biological Engineering, Northwestern University, Evanston, IL; <sup>2</sup>Department of Molecular Biosciences, Northwestern University, Evanston, IL; <sup>3</sup>Department of Chemistry, Northwestern University, Evanston, IL; and <sup>4</sup>Proteomics Center of Excellence, Northwestern University, Evanston, IL

<https://pamspublic.science.energy.gov/CCBond>

**Project Goals: The goal of this project is to develop a predictive machine learning model that will elucidate the catalytic landscape of thiamine-diphosphate dependent carboligase enzymes across the chemical space of  $\alpha$ -ketoacid substrates. This model will be used to identify enzymes for biosynthetic applications and predict potential effects of selected enzymes on the *E. coli* metabolome.**

Abstract: Increasing recognition of non-canonical enzyme activity has revealed potential problems for heterologous expression; however, understanding the potential cell burden due to promiscuous enzyme activity remains a challenge. Toward this end, our team seeks to develop cheminformatics tools that predict enzyme substrate promiscuity and predict the resulting metabolomic consequences. Specifically, we focus on the biological activity of a family of thiamine-diphosphate dependent enzymes capable of catalyzing the condensation of  $\alpha$ -ketoacids and aldehydes. Molecular fingerprints allow us to characterize the chemical space of  $\alpha$ -ketoacid substrates, from which we sample and screen for enzyme activity in high-throughput. Using an active learning approach, we plan to train a support vector machine learning model on this bioactivity data. Robust models for each carboligase enzyme will not only allow us to identify novel biosynthetic reactions but also apply predicted promiscuous activity to genome-scale models of host organisms such as *E. coli*. We plan to carry out flux balance analysis to characterize effects of side reactions and select enzymes with minimal cell burden.

*This work is supported by DOE grant DE-SC0019339.*